

Review Article

Essential Content of Software Effort Estimation using Active Learning

Divya Agarwal

Research Scholar Alagappa Chettiar College of Engineering & Technology, Karaikudi(Anna University).

I N F O

E-mail Id:

divya23@gmail.com

Orcid Id:

<https://orcid.org/0009-0002-4564-9761>

How to cite this article:

Agarwal D. Essential Content of Software Effort Estimation using Active Learning. *J Adv Res Embed Sys* 2023; 10(2): 6-14.

Date of Submission: 2023-11-11

Date of Acceptance: 2023-12-11

A B S T R A C T

Do we always need to estimate software effort (SEE) using complex methods? The objective is to characterize the core elements of SEE data, that is, the minimal number of traits and examples required to fully encapsulate the data's meaning. If the amount of important information is minimal, then: 1) the content must be brief; and 2) the value added by complex learning methods must be minimal. Technique: Our QUI Does estimate software effort (SEE) necessarily require the use of sophisticated techniques? The goal is to define the essential components of SEE data, i.e., the bare minimum of characteristics and instances needed to completely capture the meaning of the data. If there is little to no important information, then: 1) the content needs to be concise; and 2) there shouldn't be any benefit from using sophisticated learning techniques. Method The CK approach first determines the Euclidean distance between the SEE data's rows (instances) and columns (features), after which it eliminates synonyms (similar features) and outliers (far instances). Finally, it evaluates the reduced data by comparing the predictions of 1) a state-of-the-art learner (CART) using all the data, and 2) a simple learner using the reduced data. Hold-out studies are used to measure performance, which is then expressed as mean and median MRE, MAR, PRED (25), MBRE, MIBRE, or MMER. Regarding eighteen datasets, QUICK reduced the training data from 69 to 96 percent (median = 89 percent). K 14 1 closest neighbour predictions performed as well in the reduced data as did CART's predictions (using complete data). In summary, certain SEE datasets provide comparatively little essential information. Complex estimation algorithms should be simplified for such datasets as they may be unduly complex. See QUICK as an illustration of a less complex SEE strategy.

Keywords: Index Terms Software Cost Estimation, Active Learning, Analogy, K-NN

Introduction

Accurate Software Effort Estimating (SEE) is necessary for many business processes, including budgeting, project planning, iteration plans, pricing strategies, investment assessments, and bidding rounds. Expert-based techniques that provide forecasts utilizing human expertise (perhaps enhanced by process guidelines, checklists, and data) can be used to create such estimations.^{1,2} Conversely, model-based systems might use data miners to summarize current data and produce forecasts about upcoming projects.^{3,4}

Particular focus in the SEE literature is on model-based techniques for summarizing historical data. According to⁵ Jorgensen and Shepperd's SEE literature review Sixty-one percent of the selected studies deals with the introduction of new methods and how they compare to those from the past. Many corporate activities, including as budgeting, project planning, iteration plans, pricing strategies, investment assessments, and bidding rounds, require accurate Software Effort Estimation (SEE). Such estimates can be produced by expert-based techniques that generate forecasts using human expertise (perhaps strengthened by process rules, checklists, and data).^{1,2} On the other hand, data miners may be used by model-based systems to compile existing data and generate projections for future initiatives.^{3,4}

The SEE literature focuses in particular on model-based methods for historical data summarization. As stated in.⁵ Review of the SEE literature by Jorgensen and Shepperd Sixty-one percent of the chosen papers address the introduction of new techniques and their comparative analysis with historical techniques. This paper presents a novel method called QUICK for searching for NO and F O. The Euclidean distance between rows (instances) in the SEE dataset is calculated using QUICK. To find the spacing between matrix columns (features), a transposed duplicate of the matrix is utilized. Next, QUICK removes outliers (rows that are excessively distant from the rest) and synonyms (features that are extremely similar to other attributes). QUICK then makes one last use of the distance calculations to generate test instance estimates by utilizing the closest neighbour in the smaller region. The more complex the estimation processes are, the more prone to error by the operator. This is a growing problem. Shepperd et al.¹² state that the person using the data miner (rather than the dataset being investigated or the data miner being employed) is the main determinant of approach performance. This is a very worrying finding, suggesting that our sophisticated data mining methods have become so complex that they are now challenging and prone to errors. According to this research, estimate technique complexity should only be necessary if the extra advantage justifies it.

The rest of this essay is structured as follows: The symbols used in this paper are listed in Table 1. Section 2 talks

Symbol	Meaning
D	Denotes a specific data set
D'	Denotes the reduced version of D by QUICK
D^T	Transposed version of D
N	Number of instances in D
N'	Number of selected instances by QUICK from D
F	Number of independent features in D
F'	Number of independent features by QUICK from D
P	Represents a project in D
F_{cat}	Represents a feature in D
i, j	Subscripts used for enumeration
P_i, P_j	i^{th} and j^{th} projects of D , respectively
DM	An $N \times N$ distance matrix that keeps distances between all projects of D
$DM(i, j)$	The distance value between P_i and P_j
k	Number of analogies used for estimation
E_{NN}	Everyone's nearest-neighbor matrix, which shows the order of P 's neighbors w.r.t. a distance measure
$E(k)$	$E(k)[i, j]$ is 1 for $i \neq j$ and $E_{NN} \leq k$, 0 otherwise
$E(1)$	Specific case of $E(k)$, where $E(1)[i, j]$ is 1 if j is i 's closest neighbor, 0 otherwise
$E(1)[i, j]$	The cell of $E(1)$ that corresponds to i^{th} row and j^{th} column
$Pop(j)$	Popularity of j : $Pop(j) = \sum_{i=1}^n E(1)[i, j]$
n	Number of projects consecutively added to active pool, i.e. a subset of D , i.e. $n \leq N$
Δ	The difference between the best and the worst error of the last n instances.
x_i, \hat{x}_i	Actual and predicted effort values of P_i , respectively.
P_{perf_i}, P_{perf_j}	i^{th} and j^{th} performance measures
M_i, M_j	i^{th} and j^{th} estimation methods

about the related work. In Section 3, QUICK is introduced along with its application to a hypothetical instance. Our methodology is explained in Section 4. The results of the experiments are presented in Section 5. Section 7 looks at validity concerns, whereas Section 6 focuses on the findings of sanity checks on private datasets. This work is discussed from the perspective of industry practitioners in Section 8. Conclusions and a list of future studies are included in Sections 9 and 10. Table-1 Symbol used in this article

Literature Review

Active Learning

In active learning, events are arranged from most to least interesting using a heuristic (in our case, the popularity value of each row). After then, the data is examined in the order that it was sorted. Learning can be stopped early if the results from all N examples do not outperform those from a selection of M instances.

Active learning research is widely available in the machine learning literature. For instance, Dasgupta¹³ looks for assurances of generalizability in active learning. He proved that an active learning heuristic driven by greed might produce performance values that are on par with any alternative heuristic in terms of reducing the quantity of labels required.¹³ In addition, QUICK uses an original heuristic—instance popularity—to choose which examples to label first. Moreover, QUICK achieves performance comparable to supervised learners, which is similar

to Dasgupta's active learning method. Active learning performs better than supervised learning with noticeably less samples, as shown by Balcan et al.¹⁴ The results of an active learning solution are shared by QUICK., for example, comparable results to supervised learners with noticeably less sample sizes, while using a different heuristic and being accustomed to a different kind of dataset (SEE datasets) than Balcan et al.'s work. Active learning was used as an illustration of a practical application by Wallace et al.¹⁵ They recommended a citation screening strategy based on active learning.

Instance and Feature Subset Selection (FSS)

QUICK's outlier pruning can alternatively be seen as an unsupervised instance selection procedure, similar to 2.2 Feature Subset Selection (FSS) and Instance. While unsupervised algorithms operate on unlabeled data, supervised algorithms require instance labels. Prototypes are subsets of data that most accurately capture the predictive characteristics of the full dataset and are produced via instance selection techniques. A common result in instance selection is that the majority of rows in a data matrix can be removed without compromising the prediction ability of the rules learned from the remaining data. For instance, Chang's¹⁹ prototype generators examined three datasets totaling 51,415,066 instances. which, in turn, were cut down to 7, 9, and 9 percent of the initial data. For example, Li et al.⁹ use a genetic algorithm to determine which occurrences yield the best estimations. Kocaguneli et al.'s TEAK method²⁰ grouped instances before choosing clusters with minimal class label variance. QUICK and supervised instance selection techniques differ in that the former does not necessitate comprehensive costing information for every instance, while the latter demands project "labels" for every example.

The synonym pruning performed by QUICK can also be seen as an unsupervised FSS method. It is widely acknowledged showed the estimation performance can be enhanced by choosing a subset of the characteristics in the SEE dataset. Lum et al.²¹ for instance, document and capture the best practices in SEE. FSS is one of the core recommendations made by the best practices. They demonstrate how supervised and manual FSS techniques both enhance estimating performance.²¹ Principal component analysis (PCA)^{11,23} and stepwise regression (SWR)²² are two more FSS examples. To use QUICK, one only needs to be familiar with array normalization and euclidean distance calculation. However, PCA requires that the user comprehend the fundamentals of feature orthogonal transformation and correlation.²⁴ Instead of producing a subset of the individual features that a user sees on the datasets, PCA produces a new collection of—less correlated—features

(principal components), which are linear combinations of the original features. This is in contrast to QUICK.

Quick

QUICK is an active learning strategy that aids in the reduction of data interpretation difficulty by discovering the important content of SEE datasets. QUICK operates as follows:

1. Sort rows and columns according to their similarity.
2. Remove unnecessary columns (synonyms) that are very similar.
3. Remove outlier rows (outliers) that are too far apart,
4. Generate an estimate using the nearest example in the remaining data.

The sections that follow give more information on these steps.

Pruning Synonym

Synonyms are characteristics that are closely related. QUICK eliminates the following superfluous features:

Transpose the dataset matrix in step one. Depending on how the initial dataset is saved, this step may or may not be required. The rows of SEE datasets, on the other hand, often reflect historical project occurrences, whilst the columns provide the attributes that define these projects. When such a matrix is transposed, the columns represent project instances and the rows represent project characteristics. To reduce the needless effect of big numbers in the following step, columns are normalized to the 0-1 interval before transposing.

Step 2: Make a distance matrix. The corresponding distance matrix for the transposed dataset DT of F instances (DM) is a F F matrix that maintains the distances between each feature pair using the Euclidean distance function. A cell placed in the ith row and jth column (DM_{i,j}), for example, maintains the distance between the ith and jth features (diagonal cells (DM_{i,i}; I am disregarded).

Step 3: Create the ENN and E1 matrices. ENN 12i; j displays the neighbour rank of "j" in relation to I for example, if "j" is "i's" third nearest neighbour, then ENN 12i; j]14 3. The simple scenario in which I 14 j is disregarded, that is, an instance's nearest neighbour does not include itself. The following is the definition of the Ek matrix: If I j and ENN 12i; j] k, then Ek12i; j]14 1; otherwise, Ek12i; j]14 1. In synonym trimming, we aim to identify unique characteristics that do not have any nearest neighbours. To that end, we begin with k 14 1; so, E(1) distinguishes between features that have at least another nearest neighbour and those that do not. Popular features are those that show as one of the k-closest neighbours of another feature. "PopFeatj," the "popularity index" (or simply "popularity") of feature "j,"

Step 4: Using E1, compute the popularity index and select nonpopular characteristics. Nonpopular features are those with a popularity of zero, for example, PopFeat14 0.

Pruning Outliers

Outlier pruning is similar to synonym pruning, but there are some key differences: We transpose the data using synonym trimming to get the distances between “rows” (which in the transposed data are features). Then we count the popularity of each feature and remove the most popular ones (these are the features that needlessly repeated the information found in other features). We do not transpose the data before calculating the distances between rows while using outlier trimming. Then we count each row’s popularity and remove the unpopular rows (the instances that are most distant from the others). It should be noted that the dataset used to reduce outliers only comprises the previously specified characteristics. Also, from now on, the words feature and variable shall be used interchangeably.

The steps of the outlier pruning phase are as follows:

Step 1: Create a distance matrix. The related DM for a dataset D of size N is a N N matrix whose cell placed at row I and column j (DM I j) maintains the distance between the ith and jth instances of D. The diagonal cells (DMi; I am disregarded. It should be noted that current D is from the synonym trimming phase, thus it only includes the specified properties.

Step 2: Create the ENN and E1 matrices. ENN 12i; j] displays the neighbour rank of “j” in relation to “i.” comparable to the step of If “j” is “i’s” third nearest neighbour, then ENN 12i; j]14 3. Once again, the trivial case of I 14 j is omitted (nearest neighbour does not include itself).

Step 3: Calculate the popularity index

Step 4: Find stopping point and halt

Table-2: The Percentage of the Popular Instances (to Dataset Size) Useful for Prediction in a Closest Neighbor Setting

Dataset	% popular instances
kemerer	80
telecom1	78
nasa93_center_1	75
cocomo81s	73
finnish	71
cocomo81e	68
cocomo81	65
nasa93_center_2	65
nasa93_center_5	64
nasa93	63
cocomo81o	63
sdr	63
desharnaisL1	61
desharnaisL2	60
desharnaisL3	60
miyazaki94	58
desharnais	57
albrecht	54
maxwell	13

Project	Feat ₁	Feat ₂	Feat ₃	Effort
P ₁	1	2	20	3
P ₂	2	4	10	4
P ₃	3	6	40	7

Figure 1. Three projects defined by three independent features/variables and a dependent variable

	P ₁	P ₂	P ₃
Feat ₁	0.0	0.5	1
Feat ₂	0.0	0.5	1
Feat ₃	0.3	0.0	1

Figure 2. Matrix after normalizing and transposing D Example

This section provides a brief demonstration of QUICK. Assume that the example’s training set consists of three instances/projects: P1, P2, and P3. Assume that each of these projects has one dependent feature and three independent features. Our dataset would resemble the one depicted in Fig. 1.

Synonym Pruning

Step 1: Transpose the dataset matrix in step one. The resultant matrix would appear like Fig. 2 after normalising to the 0-1 interval and transposing our dataset.

Step 2: Make DMs. The distance between features is maintained by the DM. The DM in Fig. 3 is calculated using Fig. 2.

Step 3: Create the ENN and E1 matrices. According to the DM in Fig. 3, the resultant ENN 12i; j] in Fig. 4 indicates the feature neighbour rankings. Using ENN, we compute the E1 matrix (Fig. 5), which detects features that have at least one additional nearest neighbour.

Step 4: Using E1, compute the popularity index and select nonpopular characteristics. The popularity of a feature is the sum of E1’s columns (see the summation in Fig. 5). Nonpopular characteristics are those that have no popularity. We only used Feat3 in this toy example since it is the only column with zero popularity.

Outlier Removal and Estimation

QUICK continues with only the selected features during this phase. Our dataset now looks like the one in Fig. 6.

Step 1. The first QUICK step in this phase is to construct the DM. Because projects are described by a single attribute, Feat3, the Euclidean distance between two projects is the difference in normalised Feat3 values. The resulting DM is shown in Fig. 7.

Step 2. The second step is to generate the ENN matrix from the DM. We traverse the DM row by row, labelling the instances according to their distance order: the closest

neighbour is labelled 1, the second closest neighbour is labelled 2, and so on.

Step 3: Determine the labelling order and calculate the popularity index based on ENN. Remember from the previous section that E1 is generated by ENN: E112i; j]14 1.

Step 4: Determine your stopping point. Figure 10 depicts the change in the active pool for the toy example. In practise, we only move from Roundi to Roundi1 if the stopping rules are not met.

	Feat ₁	Feat ₂	Feat ₃
Feat ₁	na	0.0	0.6
Feat ₂	0.0	na	0.6
Feat ₃	0.6	0.6	na

Figure 3. DM for feature

	Feat ₁	Feat ₂	Feat ₃
Feat ₁	na	1	2
Feat ₂	1	na	2
Feat ₃	1	1	na

Figure 4. The ENN matrix for features, resulting from the DM of Figure 3. Diagonal cells are ignored

	$Feat_1$	$Feat_2$	$Feat_3$
$Feat_1$	0	1	0
$Feat_2$	1	0	0
$Feat_3$	1	1	0
$+$			
$Popularity :$	2	1	0

Figure 5. Popularity of the features. Popularity is the sum of the E I matrix columns

Project	Feat ₃	Effort
P ₁	20	3
P ₂	10	4
P ₃	40	7

Figure 6. Three projects defined by Feat3 (say, KLOC) and effort (say, in staff months)

	P ₁	P ₂	P ₃
P ₁	0	0.3	0.7
P ₂	0.3	0	1
P ₃	0.7	1	0

Figure 7. The DM of the projects P1, P2, and P3

	P ₁	P ₂	P ₃
P ₁	na	1	2
P ₂	1	na	2
P ₃	1	2	na

Figure 8. The ENN matrix resulting from the DM of Fig. 7. Diagonal cells are ignored.

		P_1	P_2	P_3
	P_1	0	1	0
	P_2	1	0	0
	P_3	1	0	0
+	Popularity :	2	1	0

Figure 9. Popularity is the sum of EglŞ's columns

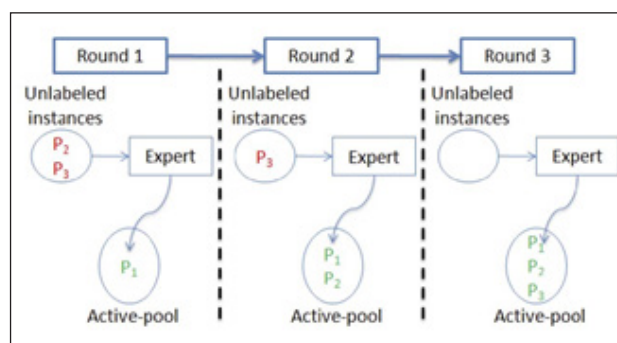


Figure 10. The change in the active pool for the toy example. In an actual setting, the transition between Roundi to Roundi1 is governed by the stopping rules.

Methodology

Datasets

Our study makes use of a total of 18 datasets (listed in Table 3), all of which are COCOMO datasets (cocomo*, Nasa*) collected using the COCOMO approach.²⁵ Three of these datasets (nasa93 centre 1, nasa93 centre 2, and nasa93 centre 5) are from NASA development centres located throughout the United States. Three other datasets are primarily from aerospace companies in southern California (cocomo81e, coco- mo81o, cocomo81s). An important point to mention here is the handling of nominal values in datasets collected using the COCOMO method. The nominal values in the COCO- MO datasets can be "low," "high," "very-high," and so on, and these values have corresponding numeric values, as explained by Boehm.²⁵ We converted nominal values to numeric equivalents in our paper. The desharnais dataset, which contains software projects from Canada, is another widely used dataset in SEE. It is gathered using the function points method. Be mindful that three projects in the Desharnais dataset have missing values. The two most common methods for dealing with missing values are 1) deleting projects with missing data, or 2) using imputation²⁴ to extrapolate the missing values from finished projects. We chose the second choice and imputed mean values in place of the missing values. Data from current initiatives of several Turkish software businesses are included in SDR. SDR is one of the newest datasets utilised in this study and is gathered by Softlab, the Bogazici University Software Engineering Research Laboratory.²⁶ Details on the IBM initiatives that make up the

albrecht dataset are provided in.²⁷ The Finnish dataset was compiled by a single person and includes 40 projects from several firms. We utilise 38 examples since the two projects with missing information are not included in our analysis.²⁸ has further information. Kemerer is a 15-instance dataset that is rather modest; further information is available in.²⁹ Another recent dataset from the financial industry, called maxwell, is made up of Finnish banking software projects and dates from the late 1990s to the early 2000s. Information is provided in.³⁰ projects created in COBOL may be found in Miyazaki. To learn more, check.³¹

Algorithm

We employ the closest neighbour and CART approaches in this investigation. Here's how we defend our decision to use SEE methods: Numerous studies come to the conclusion that the closest neighbour and CART approaches are effective SEE comparison algorithms. Walker den and Jeffrey³² support CART as a cutting-edge SEE technique. The best effort estimates approaches have been discussed in two recent works that were published in the IEEE Transactions on Software Engineering^{11,33} According to these studies, current approaches for evaluating new effort estimation methods, such as CART's regression tree generation, may be more than adequate. For instance, Dejaeger et al.³³ found little evidence that learners that are more complex than CART provide appreciable value addition. Our own findings support the findings of Walker den et al., who examined 90 effort estimators and built ensemble techniques using all conceivable permutations of 10 preprocessors and 9 learners.¹¹ Normalization, different discretises, and feature selectors served as the preprocessors. The learners comprised neural nets, linear and SWR, CART, and k-NN (with k 14 1, k 14 5). The ranking of the estimators changed across several datasets and accuracy estimators, as would have been expected by Shepperd and Kadoda.³⁴ However, we discovered a small set of 13 estimators that consistently outperformed the others (all variations of CART and k-NN assisted by preprocessors). These preprocessors and learners made up these techniques in some way.

In¹¹ two types of learners were examined: an iterative dichotomizer (CART), and an instance-based learner (ABEO-kNN). The characteristic that best divides the data so that the variance of each division is minimised is found via iterative dichotomizers like CART. After that, the process repeats on each division. In order to get the estimate, the cost information for the instances in the leaf nodes is finally averaged.

Our term for a very fundamental kind of ABE that we obtained from many ABE experiments^{9,35,36} is ABEO-kNN.

ABEO- kNN measures the distances between test and training examples using a euclidean distance function after independent variables have been initially normalised to the 0–1 range. The two preprocessors, the learners, and the two algorithms—log&ABEO-1NN and norm&CART—are mixed. We'll employ two separate iterations of log&ABEO-1NN: one that operates on the so-called active pool, which is a pool that only contains examples that have been labelled by QUICK, and another that operates on a training set that comprises every instance that has been labelled. We shall refer to the former as "activeNN" and the latter as "passiveNN" for convenience. Keep in mind that our QUICK algorithm is activeNN. The learner's name (CART) and the algorithm name (norm&CART) will now be used interchangeably because there is only one CART-based algorithm (norm&CART). Note that just two of the four potential combinations produced by the two pre-processors and two learners are used. In another article³⁷ we used seven distinct error metrics to compare all four potential combinations (along with a total of 90 different approaches) on a total of 20 datasets. We select log&ABEO-1NN and norm&CART since the comparison revealed that they perform better than the other two combinations.

Experimental Design

Our test setup consists of three components:

1. Produce baseline findings. On the complete training set, apply CART and passi-veNN.
2. Produce the outcomes of active learning. RUN FAST.
3. Compare the outcomes of QUICK to the results from the baseline.

1. Applying CART and passiveNN to the full training set will produce baseline results. On the whole training set, the algorithms are applied, and their estimates are recorded. We employ ten-way cross validation.

- a. Randomize the dataset's instance order.
- b. Create 10 bins from the dataset.
- c. Use the remaining bins as the training set and one bin at a time as the test set.
- d. Rerun the previous step using each of the 10 bins as a test set individually.

2. Run QUICK to produce the active learning outcomes. Each iteration starts with the features being chosen, then training instances are added to the active pool in the order of popularity. Only instances in the current pool may be used by QUICK, while training instances outside the pool are regarded as unlabelled. Ten-way cross validation is used to construct the train and test sets.

```

if Wilcoxon( $Perf_i$ ,  $Perf_j$ , 95) says there is no statistical difference. then
     $tie_i = tie_i + 1$ ;
     $tie_j = tie_j + 1$ ;
else
    if better(  $Perf_i$ ,  $Perf_j$ ) 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

```

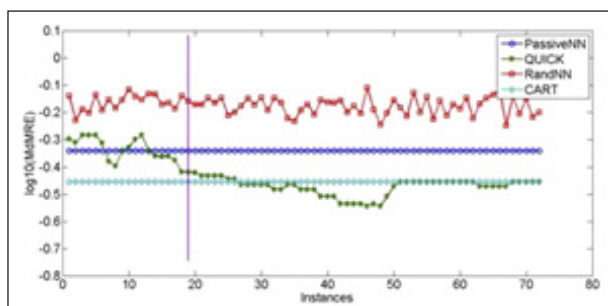


Figure 12. Sample plot of a representative (in this case, desharnais) dataset showing the stopping point (the line parallel to the y-axis at x 19) and MdmRE values (logged with base 10 for easier visualization). Note that QUICK uses only the four selected features of desharnis dataset, whereas other methods use all 10 of the features

Results

Estimation Performance

A sample plot for a representative dataset is provided in Fig. 12. (shown is the desharnais dataset). It is the outcome of labelling NO occurrences in decreasing popularity after using QUICK's synonym trimming (four characteristics were used for desharnais). Here is the reading from Fig. 12:

The recorded MdmRE error measurements are displayed on the Y axis. The performance improves with decreasing value.

The line with star dots displays the mistake that resulted from using the labels 1... I — 1 to estimate the ith occurrence.

The horizontal lines display the inaccuracies that were made when estimations were made utilising all the available data (either from CART or passiveNN).

The vertical line indicates the location at which QUICK suggested labelling can end (i.e., NO).

The square-dotted line represents randNN, which is produced by choosing any instance at random from the training set and using its effort value as the estimate.

Fig. 12 highlights three key findings: 1) the decrease in the number of occurrences necessary, 2) the 3) The estimation inaccuracy, together with a reduction in the number of characteristics. QUICK achieves the same low error rates as CART and passiveNN, both of which employ the complete dataset, using a portion of the occurrences and characteristics of the original dataset.

Comparison QUICK versus CART

The PRED (25) difference between CART (using all the data) and QUICK is displayed in Fig. 13. (using just a subset of the data). PRED (25) of CART less PRED (25) of QUICK equals the difference. Therefore, a negative value means that PRED (25) estimations from QUICK are superior than those from CART. Starting with the value of —35, the left-hand side demonstrates that QUICK performs better than CART, but up to the value of 35, the right-hand side demonstrates that CART outperformed QUICK in certain instances.

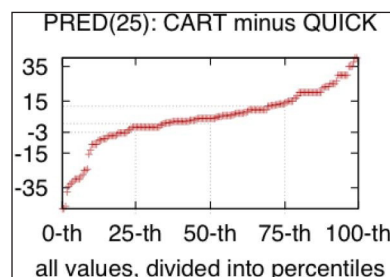


Figure 13. CART minus QUICK values for PRED(25) in 18 datasets. Negative values mean that QUICK outperformed CART. The median is only 2 percent with a small interquartile range (75th 25th percentile) of 15 percent.

According to Fig. 13, the performance of CART and QUICK are quite similar at the median point, with the 50th percentile corresponding to a PRED (25) value of roughly 2. Also take note that the 75th percentile is less than 15, indicating that the difference between CART and QUICK is not very noticeable when one is superior to the other. Our findings indicate that employing all projects and features will only offer a little amount of value. Estimates as accurate as those produced by more complicated learners like CART can be obtained by doing a QUICK analysis on a tiny portion of the data.

Detailed Statistical Analysis

QUICK is compared to passiveNN and CART using seven assessment criteria in Tables 5 and 6. In these tables, smaller values are preferable. For six of the metrics, the term "loss" denotes larger error levels. For PRED (25), however, "loss" refers to lower values. Each table's last column totals the

method's losses in the corresponding row. Table 5's final column may be sorted to reveal that the loss figures are extremely similar:

QUICK: 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 3; 6

passiveNN: 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 2; 3; 6; 6; 7:

Table 5's grey rows represent the datasets where QUICK fails to outperform most error measures (4 times out of 7 error measures, or more). The most important finding is that, in comparison to a thorough analysis of all projects, a QUICK analysis loses less often (just 1 grey row) when employing closest neighbour algorithms. As seen in Fig. 13, QUICK performs similarly to CART. This is also evident in Table 6's last column, which totals the number of losses. The datasets where QUICK loses to CART most frequently (4 or more times out of 7) are represented by the four grey rows in Table 6. Only 4=18 14 22% of the datasets benefit more from a complete CART analysis than a QUICK partial analysis of a small subset of the data. The sorted last column in Table 6 is

QUICK :0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 1; 3; 6; 7; 7; 7

CART :0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0; 0:

Dataset	Method	MMRE	MAE	Pred(25)	MMRE	MBRE	MBRE	MMER	SUM OF LOSSES
cocomo01	CART	0	0	0	0	0	0	0	0
	QUICK	1	1	1	1	1	1	1	7
cocomo01a	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
cocomo01b	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
cocomo01c	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
desfarmis	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
desfarmis_L1	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
desfarmis_L2	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
desfarmis_L3	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	1	1	1	3
nasa01	CART	0	0	0	0	0	0	0	0
	QUICK	0	1	0	0	0	0	0	1
nasa01_center_1	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
nasa01_center_2	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
nasa01_center_3	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
sair	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
allsrcfit	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
hannush	CART	0	0	0	0	0	0	0	0
	QUICK	1	1	1	1	1	1	1	7
kennecor	CART	0	0	0	0	0	0	0	0
	QUICK	0	0	0	0	0	0	0	0
maxxviii	CART	0	0	0	0	0	0	0	0
	QUICK	1	1	1	1	1	1	1	7
mryazaki94	CART	0	0	0	0	0	0	0	0
	QUICK	1	0	1	1	1	1	1	6

8. A. Corazza, S. Di Martino, F. Ferrucci, C. Gravino, F. Sarro, and
9. E. Mendes, "How Effective Is Tabu Search to Configure Support Vector Regression for Effort Estimation?" Proc. Sixth Int'l Conf. Predictive Models in Software Eng., p. 4, 2010.
10. Y. Li, M. Xie, and T. Goh, "A Study of Project Selection and Feature Weighting for Analogy Based Software Cost Estimation,"
11. J. Systems and Software, vol. 82, no. 2, pp. 241-252, 2009.
12. T. Menzies, Z. Chen, J. Hihn, and K. Lum, "Selecting Best Practices for Effort Estimation," IEEE Trans. Software Eng., vol. 32, no. 11, pp. 883-895, Nov. 2006.
13. E. Kocaguneli, T. Menzies, and J. Keung, "On the Value of Ensemble Effort Estimation," IEEE Trans. Software Eng., vol. 38, no. 6, pp. 1403-1416, Nov./Dec. 2012.
14. M. Shepperd, "It Doesn't Matter What You Do but Does Matter Who Does It!" Proc. CREST Open Workshop, Oct. 2011.
15. S. Dasgupta, "Analysis of a Greedy Active Learning Strategy,"
16. Proc. Neural Information Processing Systems, vol. 17, 2005.
17. M.-F. Balcan, A. Beygelzimer, and J. Langford, "Agnostic Active Learning," Proc. 23rd Int'l Conf. Machine Learning, pp. 65-72, 2006.
18. B. Wallace, K. Small, C. Brodley, and T. Trikalinos, "Active Learning for Biomedical Citation Screening," Proc. 16th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 173-182, 2010.
19. J.F. Bowring, J.M. Rehg, and M.J. Harrold, "Active Learning for Automatic Classification of Software Behavior," ACM SIGSOFT Software Eng. Notes, vol. 29, pp. 195-205, July 2004.
20. T. Xie and D. Notkin, "Mutually Enhancing Test Generation and Specification Inference," Proc. Int'l Workshop Formal Approaches to Software Testing, pp. 1100-1101, 2004.
21. A. Hassan and T. Xie, "Software Intelligence: The Future of Mining Software Engineering Data," Proc. FSE/SDP Workshop Future of Software Eng. Research, pp. 161-166, 2010.
22. C.L. Chang, "Finding Prototypes for Nearest Neighbor Classifiers," IEEE Trans. Computers, vol. 23, no. 11, pp. 1179-1185, Nov. 1974.
23. E. Kocaguneli, T. Menzies, A. Bener, and J.W. Keung, "Exploiting the Essential Assumptions of Analogy-Based Effort Estimation," IEEE Trans. Software Eng., vol. 38, no. 2, pp. 425-438, Mar./Apr. 2012.