

Local vs. Global Lessons for Defect Prediction and Effort Estimation

Tim Menzies, *Member, IEEE*, Andrew Butcher, David Cok, *Member, IEEE*,
Andrian Marcus, *Member, IEEE*, Lucas Layman, Forrest Shull, *Member, IEEE*,
Burak Turhan, *Member, IEEE*, Thomas Zimmermann

Abstract—Existing research is unclear on how to generate lessons learned for defect prediction and effort estimation. Should we seek lessons that are global to multiple projects, or just local to particular projects? This paper aims to comparatively evaluate local vs. global lessons learned for effort estimation and defect prediction. We applied automated clustering tools to effort and defect data sets from the PROMISE repository. Rule learners generated lessons learned from all the data, from local projects, or just from each cluster. The results indicate that the lessons learned after combining small parts of different data sources (i.e., the clusters) were superior to either generalizations formed over all the data or local lessons formed from particular projects. We conclude that when researchers attempt to draw lessons from some historical data source, they should (a) ignore any existing local divisions into multiple sources; (b) cluster across all available data; then (c) restrict the learning of lessons to the clusters from other sources that are nearest to the test data.

Index Terms—Data mining, clustering, defect prediction, effort estimation

1 INTRODUCTION

PROCESS and product data are used in software engineering (SE) to support a variety of tasks, such as, defect prediction, effort estimation, refactoring of source code, determination of the social networks of programmers, learning the expected interface between modules, etc. Two questions are at the center of much of the research in the field: (a) what data (e.g., which metrics) are best suited to support specific tasks?; and (b) what is the best way to reason about SE data? This paper explores the latter question, in the context of:

- **Software effort reduction:** finding rules for reducing a project's development time;
- **Software defect reduction:** finding rules for reducing a project's defect count.

Our focus in this paper is not on what data (e.g., process or product data) are used for building models for defect prediction or effort estimation, but rather on the source of the data and implicitly the applicability of the lessons derived by the models. Software *data* come from some *sources* (e.g., different companies for effort data or different projects for defect data). These data show the defects or effort associated with

examples (e.g., projects for effort data or classes for defect data) from that source.

How should we reason about these data? On this point, the literature is contradictory. Some existing work argues that *data from multiple sources* can generate rules that apply in any context (i.e., project or company) [1], [2], [3], [4]. We call these *global lessons* (or rules).

On the other hand, other papers indicate that the best lessons are learned from *within one source* (rather than across all sources), which implies that these lessons are only useful in their context [5], [6], [7]. We call these *local lessons*.

When project managers want to make changes in their projects in order to minimize the development effort or the rate of defects, they are faced with two options: (a) make changes based on *global lessons* available from existing data; or (b) mine data about the current project and infer *local lessons*. The dilemma of the manager is obvious. In the first case, expensive changes (based on the global lessons) may be undertaken without reaching the desired goal, if the global lessons prove to be wrong for this context. In the second case, an upfront investment is needed to collect and analyze data and to generate the local lessons, which may be unnecessary if the global lessons apply to this context.

This paper addresses this dilemma and reports on experiments where:

- Data from different sources are combined;
- Within that combination, automatic tools find clusters of related examples. Note that clusters may contain examples from multiple sources.

T. Menzies and A. Butcher are with the Lane Department of Computer Science and Electrical Engineering, West Virginia University, Morgantown, WV; tim@menzies.us, abutcher@froglegs.com.

D. Cok is with GrammaTech, Ithaca, New York; dcok@grammatech.com.

A. Marcus is with Wayne State University; amarcus@wayne.edu.

L. Layman and F. Shull are with the Fraunhofer Institute, University of Maryland, USA; llayman@fc-md.umd.edu, fshull@fc-md.umd.edu.

B. Turhan is with the Department of Information Processing Science, University of Oulu, Oulu, Finland; turhanb@computer.org.

T. Zimmermann is with Microsoft Research, Redmond, USA; tzimmer@microsoft.com.

- Data mining is then used to learn lessons from the examples in each cluster.
- The generated lessons are compared.

Before this experiment, our previous work indicated that local lessons seem to be best for defect prediction and effort estimation [8], [9]. The surprising result of the experiment presented in this paper is that the best lessons for a project from one source come from *neighboring clusters with data from nearby sources, but not inside, that source*. We will call such lessons *neighbor lessons*.

We conclude that when project managers attempt to draw lessons (i.e., about defect or effort reduction) from some historical data source, they should (a) ignore any existing local divisions into multiple sources; (b) cluster across all available data; then (c) restrict the learning of lessons to the clusters that are nearest to the test data (regardless of the source of the data). While the paper is focusing on defect prediction and effort estimation data, we believe that our conclusions can translate to other types of data, related to different SE tasks.

This paper extends a prior publication [10] in the following way:

- That paper only explored four data sets. Here, we explore over twice that number of data sets.
- That prior publication only explored local vs. global lessons and found evidence that supported local lessons over the global ones. In this paper we offer new experiments in §4 (which have not appeared previously) that explore clusters from multiple sources.
- This paper's literature review is more extensive (see Figure 2).

The rest of this paper is structured as follows. First, we explore in §2 the literature on defect prediction and effort estimation. This highlights the lack of stable conclusions on what is the best way to create predictors (i.e., how to best learn lessons from the existing data). §4 presents two experiments that compare rules for defect prediction and effort estimation generated from global data, from strict local data, and from clusters, whereas §3 discusses the details of the clustering and rule learning algorithms used in these experiments. §5 discusses the external validity of the conclusions resulting from the two experiments.

2 LITERATURE REVIEW

THE main goal of this section is not to provide a generic review of defect prediction and effort estimation work, but to highlight work that documents contradictory results, which make it difficult to generalize solutions for defect prediction or effort estimation.

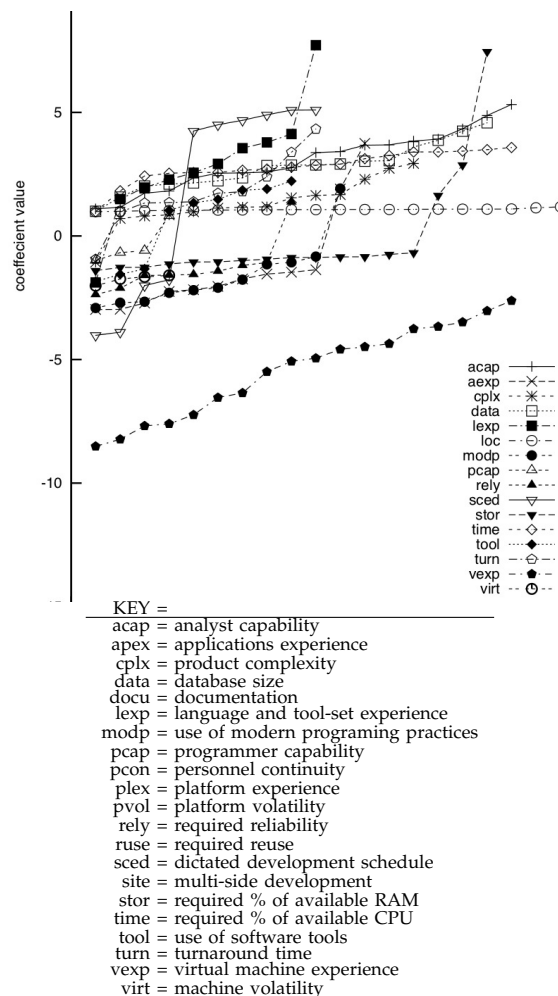


Fig. 1. Sorted β_i values from local calibration on 20*(66%) samples of NASA93 data (From [11]). Coefficients learned using Boehm's recommended methods [12]. A greedy backward selection removed attributes with no impact on estimates (so some attributes have less than 20 results).

2.1 Effort Estimation

Conclusion instability in effort estimation may be a fundamental property of the data sets we are exploring [11]. For example, Figure 1 tests the stability of Boehm's COCOMO software development effort estimation model [11]. In that analysis, 20 times, we learned $effort = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$ using a random $\frac{2}{3} rds$ sample from 93 NASA projects (this subset size was chosen to be similar in size to the 61 projects used to find the original COCOMO coefficients). As shown in Figure 1, only the coefficient on lines of code (*loc*) was stable. The observed ranges on the other β_i coefficients were very large (e.g. $-8 \leq vexp \leq -3.5$). In fact, the signs of five coefficients even changed from positive to negative (see *stor*, *aexp*, *modp*, *cplx*, *sced*). This coefficient instability is particularly troubling since we know of NASA project managers who have made acquisition decisions worth tens of millions of

ref	cbo	rfc	lcom	dit	noc	wmc	#projects	size	type
[13]	+	+	+	-	-	+	6	95-201 classes	6 versions of rhino (java)
[14]	+	+	+	-	-	+	12	86 classes (3-12kloc)	student
[15]	+	+	-	-	-	-	1	1700 classes (110kloc)	commercial telecom
[16]	+	+	-	+	+	+	8	113 classes	student
[17]	+	+	-	+	+	+	8	114 classes	student
[18]	+	+	+	+	-	-	1	83 classes	commercial: lalo (c++)
[19]				+	+	-	1	32 classes	commercial: telecom c++
[20]				+	-	-	1	42-69 classes	commercial java word proc.
[21]	+	-	-	-	-	-	1	85 classes	telecom c++
[22]	-	+	-	-	-	+	3	92 classes	3 c++ subsystems,commercial
[23]	+	+	+	-	+	+	1	123 classes (34kloc)	java commercial
[24]	+			+		+	1	706 classes	commercial c++ and java
[25]	+	+	+	-	+	+	1	145 classes	kc1-nasa
[26]	+	+	+	+	-	+	1	3677 classes	open source:mozilla
[27]	+	+	+			+	1	?	java (sap) commercial
[28]	+	+	+	+	+	+	3	?	eclipse 2.0, 2.1, 3.0
[29]	-	+	+	-	-	+	8	113 classes	student
[30]		+	+	+	+		2	64 classes	sales and cd-selection system
[31]		-	-	-	-	-	1	3344 modules (2mloc)	commercial telecom c++
[32]	+	+	+	-	-	+	5	395 classes	commercial telecom c++
[33]	+	+	-	-	-	+	1	1412 classes	open source:jdt
[34]	+	+	-	-	-	+	2	9713 classes	eclipse 2.0, 2.1
[35]	+	+	-	-	-	+	1	145 classes	kc1-nasa
[36]				+	-	-	1	145 classes	commercial java xml editor
[37]	-	-	-	-	-	-	1	174 classes	commercial telecom c++
[38]	-					-	0	50 classes	student
[39]	+	+	-	-	-	+	1	145 classes	kc1-nasa
[40]		+		+	+		2	294 classes	commercial c++
total +	18	20	11	11	8	17			
total -	4	3	7	14	16	4			

KEY:	Strong consensus (over 2/3rds)
	Some consensus (less than 2/3rds)
	Weak consensus (about half)
	No consensus

Total percents: "*" denotes majority conclusion in each column						
+	* 64%	* 71%	* 39%	39%	29%	* 61%
-	14%	11%	25%	* 50%	* 57%	14%

Fig. 2. Contradictory conclusions from OO-metrics studies for defect prediction. Studies report significant (“+”) or irrelevant (“-”) metrics verified by univariate prediction models. Blank entries indicate that the corresponding metric is not evaluated in that particular study. Colors comment on the most frequent conclusion of each column. CBO= coupling between objects; RFC= response for class (#methods executed by arriving messages); LCOM= lack of cohesion (pairs of methods referencing one instance variable, different definitions of LCOM are aggregated); NOC= number of children (immediate subclasses); WMC= #methods per class.

dollars based on the COCOMO coefficients (i.e., they decided to acquire the technologies that had most impact on the variables with largest coefficients).

Other papers also report contradictory findings about effort estimation. Jørgensen [41] reviewed 15 studies comparing model-based to expert-based methods. Five of those studies favored expert-based methods, five found no difference, and five favored model-based methods. Kitchenham et al. [42] reviewed studies that checked if data imported from other organizations were as useful as local data (for the purposes of building effort models). From a total of seven studies, three found that models from other organizations were not significantly worse than those based on local data, while four found that they were significantly different (and worse). MacDonell and Shepperd [43] also performed a review on the value of local vs. global effort estimation models through a replication of [42]. From a total of 10 studies, two were found to be inconclusive, three supported global models, and five supported local models. Similarly, Mair and Shepperd [44] compared regression to analogy methods for effort estimation and found conflicting evidence.

From a total of 20 empirical studies, seven favored regression, four were indifferent, and nine favored analogy.

2.2 Defect Prediction

In the area of defect prediction, there are also many contradictory findings. For example, Zimmermann et al. [45] learned defect predictors from 622 pairs of projects $\langle project_1, project_2 \rangle$. In only 4% of pairs, the defect predictors learned in $project_1$ worked in $project_2$. Similar findings (of contradictory conclusions) concern the OO metrics as well. Figure 2 lists 28 studies that offer contradictory conclusions regarding the effectiveness of OO metrics for predicting defects (exception: response for class is often a useful indicator of defects). To create Figure 2, we:

- 1) Used our domain knowledge to pick three high-impact seed articles [19], [20], [21];
- 2) Used Google Scholar¹ to find 500+ relevant studies that cited any of those seed articles;

1. <http://scholar.google.com>

- 3) Removed false positives by scanning titles and abstracts. This reduced the 500+ articles to 86;
- 4) Applied the following relevancy rule to reduce the 86 studies to 28: reject all papers that do not offer a univariate predictive analysis for the validation of the metric(s) under investigation.
- 5) Checked the literature reviews of important papers in this field [38], [46], for papers not in our sample.

For the manager of a software project Figure 2 is particularly troubling. Each study makes a clear, but usually different, conclusion. Hence, it is difficult for a manager to make clear decision about, for example, the merits of a proposed coding standard where maximum depth of inheritance is required to be less than some expert-specified threshold.

As to the root cause of the instabilities of Figure 2, we offer the following conjecture. We showed above in §2.1 that models learned from different regions within effort data can have very different properties. If defect data was as varied as effort data, then we would naturally expect that different samples of different projects would yield different models (e.g., as seen in Figure 2) due to *dataset shift* [47].

If this conjecture is correct, then we would expect that clusters within the data should produce different models. This paper is a test of that conjecture. In short, we show that different regions of the data generate different models. Further, the models built from specialized regions within the data set *perform better than those learned across all data*.

	dimension	attribute	project1	project2	project3
independent	1	afp	1587	260	152
	2	input	774	9	5
	3	output	260	4	3
	4	enquiry	340	3	2
	5	file	128	193	42
	6	interface	0	41	35
	7	added	1502	51	16
	8	changed	0	138	0
	9	deleted	0	61	0
	10	pdr_afp	4.7	16	4.4
	11	pdr_ufp	5	16.6	4.1
	12	npdr_afp	4.7	16	4.4
	13	npdu_ufp	5	16.6	4.1
	14	resource	4	2	1
	15	dev.type	0	0	0
	16	duration	4	17	9
dependent	1	effort	7490	5150	668

Fig. 3. Example data. Three examples from the CHINA effort estimation data set.

3 INSIDE THE LEARNERS

BEFORE we describe our experiments, this section reviews the technical details of the clustering algorithm and the rule learner used in these experiments. Note that our algorithms are agnostic with respect to the semantics of their input data, which means that we use the same algorithms for defect data and effort data.

In order to better relate the algorithms to our context (i.e., defect prediction and effort estimation) and to better understand the experiments of the subsequent sections, Figure 3 and Figure 4 provide examples of instances from one of the actual effort estimation data sets used in the experiments. In this data set, each project is a point in a 16-dimensional space of “independent” attributes. Each point also has one “dependent” attribute (in Figure 3 it is the development effort associated with each project). The attribute names in this figure concern function points of these projects and are defined in Figure 5.

Figure 4 shows the same data as Figure 3, but each value has been categorized into “hi” or “lo”, depending on whether it falls above or below the mean value for each row. This figure classifies our projects into two “hi” effort projects (project1 and project2) and one “lo” effort project (project3).

Contrast sets list the differences between classes. To find a contrast set, we look for an attribute where (1) all the values from one class are similar but (2) those values differ in different classes. The only such contrast set in Figure 4 is the row marked in gray and denotes the function points of internal logical files. Using this contrast set, we would say that “lo” file function points is the factor that most selects for low effort projects.

While the rules generated with contrast sets are simple, they are quite powerful and our choice is not

	dimension	attribute	project1	project2	project3
independent	1	afp	hi	lo	lo
	2	input	hi	lo	lo
	3	output	hi	lo	lo
	4	enquiry	hi	lo	lo
	5	file	hi	hi	lo
	6	interface	lo	hi	hi
	7	added	hi	lo	lo
	8	changed	lo	hi	lo
	9	deleted	lo	hi	lo
	10	pdr_afp	lo	hi	lo
	11	pdr_ufp	lo	hi	lo
	12	npdr_afp	lo	hi	lo
	13	npdu_ufp	lo	hi	lo
	14	resource	hi	lo	lo
	15	dev.type	lo	lo	lo
	16	duration	lo	hi	lo
dependent	1	effort	hi	hi	lo

Fig. 4. Example data of Figure 3, all data categorized as “hi” or “lo” depending on whether it is above or below the mean value for each row, respectively.

afp	adjusted function points	adjusted size by the standard value adjustment factor (vaf)
input	function points (ufp) of input	
output	function points (ufp) of external output	
enquiry	function points (ufp) of external enquiry	
file	function points (ufp) of internal logical files or entity references	
interface	function points (ufp) of external interface added	function points (ufp) of new or added functions
changed	function points (cfp) of changed functions	
deleted	function points (cfp) of deleted functions	
pdr_ufp	normalized level 1 productivity delivery rate	norm. level 1 effort (for development team) divided by functional size (unadjusted function points).
npdr_afp	normalized productivity delivery rate	normalized effort divided by functional size (unadjusted function points).
npdu_ufp	productivity delivery rate (adjusted function points)	summary work effort divided by adjusted function point count.
resource	team type	1 = development team effort (e.g., project team, project management); 2 = development team support (e.g., database administration, data administration, quality assurance); 3 = computer operations involvement (e.g., information center support, computer operators, network administration); 4 = end users or clients (e.g., user liaisons, user training time)
dev.type	development type	1= new development, 2= enhancement; 3= redevelopment.
duration	total elapsed time for the project	in calendar months.
effort	summary work effort	provides the total effort in hours recorded against the project.

Fig. 5. Function point measures used in the CHINA effort estimation data set. The last line shows the dependent attribute (total effort in hours).

accidental. We have previously conducted extensive evaluations of contrast set learning with other learners on software engineering data [48]. We found that these succinct rules out-performed more complex models generated by standard classifier or optimization algorithms [49]. Also, using stochastic sampling, contrast set learning can process very large data sets in linear time, which is an advantage when the data sets are large. Further, a major advantage of contrast set learning for software engineering is that it generates very succinct rules. For example, the single rule generated by this example is:

$$\text{if file=lo then effort=lo}$$

In [50] we argued that such brevity was important when explaining data mining results to business users and we prefer such rules (when possible) over more complex ones, which tend to be harder to understand.

In order to test this rule, we run a *selection study* where we look at the class distribution of all projects where *file=lo*. In this case, 100% of projects with *file=lo* have “lo” effort. Note that for this test to be valid, it should be conducted on data not used for learning the rule.

3.1 Clustering Data with WHERE

Our clustering algorithm, named WHERE, assumes that *the dimensions of most interest are the dimensions of greatest variability*. This assumption is shared by other researchers such as those using feature weighting based on variance [51] or principal component analysis (PCA), e.g., [4].

Matrix factorization methods like PCA take polynomial time to execute [52]. We focus in this paper

on defect and effort data, but our long term view is that these techniques can be used for other types of SE data. Hence, we adopt a more efficient solution for our tools. Faloutsos & Lin [53] offer a linear-time heuristic for generating these dimensions, which we use in our work. Given N instances, their “FASTMAP” heuristic finds the dimension of greatest variability to a line drawn between the two furthest points. These two points are found in linear time, as follows: *firstly* select any instance Z at random; *secondly* find the instance X that is furthest away from Z ; *thirdly* find the instance Y that is furthest away from X . The line \overline{XY} is an approximation to the first component found by PCA.

As shown in Figure 6.a, an orthogonal dimension to \overline{XY} can be found by declaring that the line \overline{XY} is of length c and runs from point $(0, 0)$ to $(0, c)$. Each instance now has a distance a to the origin (instance X) and distance b to the most remote point (instance Y). From the Pythagoras theorem and cosine rule, each instance is at the point $x = (a^2 + c^2 - b^2)/(2c)$ and $y = \sqrt{a^2 - x^2}$. Figure 6.a shows four quadrants defined by the median values of each dimension (\hat{x}, \hat{y}) : *NorthWest*, *NorthEast*, *SouthWest*, and *SouthEast*.

WHERE constructs Figure 6.a using a standard Euclidean distance operator, then it recurses on each quadrant. This generates a balanced tree of quadrants, stopping when a sub-quadrant has less than a *minimum number instances* (currently the square root of the total number of instances). The resulting tree of quadrants is then pruned from the leaves back to the root: leaf quadrants with *similar density* (currently, within 50%) are grouped together.

Figures 6.b shows the CHINA effort estimation data set from the PROMISE repository mapped into the

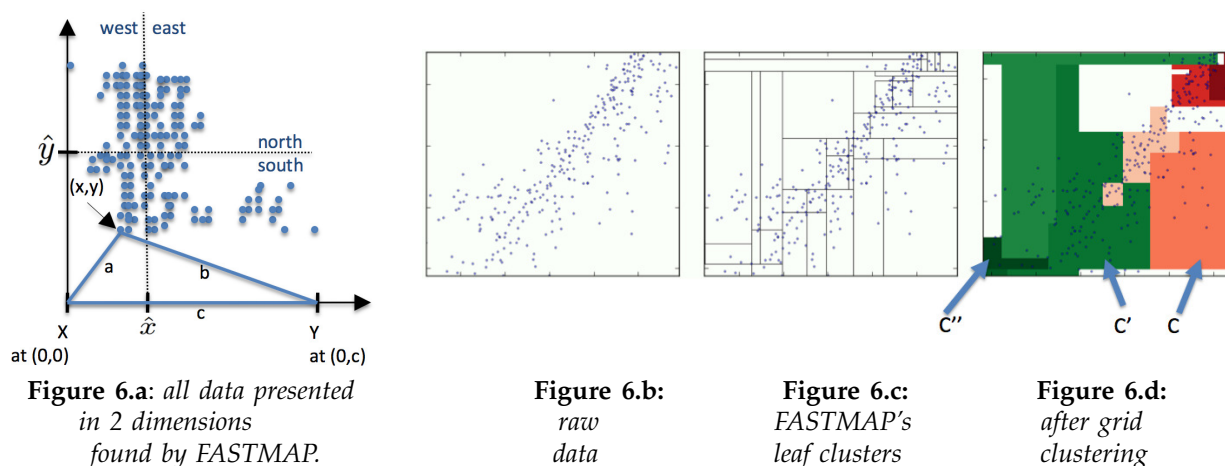


Figure 6.a: all data presented in 2 dimensions found by FASTMAP.

Figure 6.b: raw data

Figure 6.c: EASTMAP's leaf clusters

Figure 6.d: after grid clustering

Fig. 6. Each dot is a D-dimensional instance from the CHINA data set mapped into two dimensions (for an example of three such dots, see Figure 3). From that data, a new dimension is synthesized on a line between X (at the origin) and the most remote instance Y (at $0, c$)- see Figure 6.a . Each dot has distance a from the origin and b from the most remote point. The median point on the x and y axis are \hat{x} and \hat{y} , respectively. The algorithm then recurses on each quadrant to generate grids. Leaf pruning then combines the smaller clusters into the colored regions shown in Figure 6.d.

axes found by FASTMAP. Each dot describes one project using multiple independent attributes and one dependent attribute showing the development effort (in months). For an example of one of those dots, see Figure 3.

Figure 6.c shows the leaf quadrants found by WHERE's recursive exploration of the NorthWest, NorthEast, SouthWest, and SouthEast quadrants.

Figure 6.d shows the results of leaf pruning: those clusters are colored to show the median intra-cluster development effort (dark red = highest effort while dark green = lowest effort).

Now consider the three clusters labeled C, C', C'' . Suppose a manager of a project in the orange cluster C is considering how to decrease the development effort of that project (of all the neighbors of that cluster, the green cluster C' has the lowest development effort). Accordingly, that manager would learn rules over the C' data to find treatments that convert projects of type C to C' (note that such a strategy is *not* available to the manager of projects in the dark green cluster C'' because no neighbor of C'' has a shorter development effort, so there we would advise to maintain the status quo).

3.2 Learning Rules with WHICH

A vital requirement for this work is that whatever data miner is used, it can generate rules that can be compared with the general truisms in the field. Therefore, we eschew learners that use statistical distributions and probability calculations to generate models which, even if they work successfully, are opaque to a human reader. Hence, we do not use Bayes classifiers [54] or neural networks [55]. For the

same reason, we also choose not to use learners that generate numeric combinations of project influences, e.g., linear regression, logistic regression, simulated annealing [56], model trees [57], or support vector machines. Finally, we avoid learners that produce large and hard to read theories, e.g., genetic programming algorithms [58] that can learn large and intricate rules.

For this research, we use the the WHICH contrast rule learner [59]. WHICH was informally introduced at the start of §3. More formally, we say that WHICH learns rules of the form

$$\text{if } R_x \text{ then } (\text{change} = \epsilon_1/\epsilon_0 * \text{support})$$

Here, R_x is a *treatment* containing a set of attribute value pairs a_v ; ϵ_0 is the median score for instances in the untreated population. Referring back to Figure 4, we say that the untreated population is all the test data (i.e., all of Figure 4) and the treated population are the rows found by the selection study (i.e., all the examples that do not conflict with the treatment). ϵ_1 is the median score for the population subset selected by the rule. For effort and defect prediction, the ratio ϵ_1/ϵ_0 is smaller if the treatment selects for *better* instances. For an example of such a rule, see the start of this section.

WHICH builds these rules by looping over attribute value combinations, combining those that look most promising (for the AI-literature reader, we note that WHICH is a fuzzy beam search). Continuous attributes are first discretized to a few discrete values. A stack is created containing one item for every attribute value. The items in that stack are sorted using $\epsilon_1/\epsilon_0 * \text{support}$ (where *support* is the percent of the data selected by that rule). WHICH generates new rules as follows: several times, (a) select two items

at random, favoring those with better $\epsilon_1/\epsilon_0 * support$; (b) combine the pair into a new item and score it. The new rules are then sorted back into the stack. This process repeats until no new improvements are seen at the top-of-stack. WHICH returns the rule at the top-of-stack.

4 EXPERIMENTS

WE ran a set of experiments that compare the results of defect prediction and effort estimation obtained from learning rules from global data, local data, and clusters, respectively.

The objects of our experiments are the nine data sets shown in Figure 7. Each of the nine data sets used in the experiments were scored with their project effort or number of defects. There are two data sets used for effort estimation and seven used for defect prediction. The projects have three sets of attributes:

- 1) NASACOC contains the independent attributes from Figure 1. Its dependent attribute is development effort, measured in terms of calendar months (at 152 hours per month, including development and management hours).
- 2) CHINA contains the independent attributes of Figure 5 and the dependent attribute shown in the last row of that table (summary work effort).
- 3) The other seven data sets contain the independent attributes of Figure 8 and the dependent attribute of defect counts (as seen in a post-release bug tracking system).

NASACOC and CHINA are effort estimation data sets while the other seven are defect prediction data sets. Thanks to the researchers who shared the data via the PROMISE repository [60], we were able to experiment with a diverse set of data, which remains available for future replications.

The main research question we are addressing is how to generate lessons that lead to rules for minimizing effort and defects. On one hand, we generated lessons from global data and applied them to individual project data. On the other hand, we clustered the data and for each cluster, we applied lessons from the best neighboring cluster (i.e., with better defect or effort values). We present an informal example, followed by the formal description of the treatments we chose.

Assume that we have data from two sources (e.g., effort or defect data from different companies) A and B:

- Data set A = {Xa, Ya, Za}
- Data set B = {Xb, Yb, Zb}

We combine the data {Xa, Ya, Xb, Yb, Za, Zb} and cluster it. We obtain three clusters: $C = \{Xa, Xb\}$, $C' = \{Ya, Yb\}$, and $C'' = \{Za, Zb\}$. We want to see what treatments to each cluster of data result in lowering the defect (or effort) values. The question we are

About	System	Description	Size
effort	CHINA	Function points	499 projects
	NASACOC	COC81 + NASA93	156 projects
defects	JEDIT	Text editor	306 classes
	LUCENE	Text search engine	428 classes
	SYNAPSE	Enterprise service bus	256 classes
	TOMCAT	Apache servlet container	858 classes
	VELOCITY	Template language engine	229 classes
	XALAN	XSLT processor	875 classes
	XERCES	XML processor	329 classes

Fig. 7. Data from <http://promisedata.googlecode.com/> for this study.

addressing is how to infer the rules (using WHICH). We have the following options:

- 1) *Global learning*. For cluster C, WHICH learns the rules from all the data (minus what is in the cluster - {Ya, Yb, Za, Zb}) and test it on the data in C, i.e., {Xa, Xb}.
- 2) *Cluster learning*. For cluster C, WHICH learns the rules on the data from cluster C', i.e., {Ya, Yb}, (we assume here that C' is the neighbor cluster to C with the best defect/effort values) and we test it on the data in C, i.e., {Xa, Xb}. We further refine the cluster learning:
 - a) *Neighbor learning (cross)*. For the data in C from one source (e.g., Xa), WHICH learns on the data from C' that does not come from the same source (i.e., Yb in this example). Note that this is an overfitting-avoidance strategy since it somewhat increases the difference of the training data from the test data. More details on this issue are provided later in this section.
 - b) *Local learning (within)*. For the data in C from one source (e.g., Xa), WHICH learns on the data from C' that also comes from the same source (i.e., Ya in this example).

We conducted the experiment in two stages: one for steps 1 and 2 above and one for steps 2.a and 2.b. Before we describe the experiments more formally, we need to introduce some notation.

- Let *all* refer to all examples in C.
- Let *treated* refer to the examples of C selected by the rules predicting for lower class values ($treated \subseteq all$). In the case of the defect data sets, the class is “number of defects” while in the case of the effort data sets, the class is “development effort”. Note that in both cases, we wish to *minimize* the class value.
- Let $D_{treated}$ be the distribution of class variables seen in *treated*.
- Let D_{all} be the distribution of the class variables seen in *all* and *max* be the maximum value of that distribution (i.e., the worst case result). Normalize all values in D_{all} by expressing them as a percent of the *max* worst case. For example, in $D_{treated}$, a value of 50 would denote a class value

amc	average method complexity	e.g. number of JAVA byte codes
avg_cc	average McCabe	average McCabe's cyclomatic complexity seen in class
ca	afferent couplings	how many other classes use the specific class.
cam	cohesion amongst classes	summation of number of different types of method parameters in every method divided by a multiplication of number of different method parameter types in whole class and number of methods.
cbm	coupling between methods	total number of new/redefined methods to which all the inherited methods are coupled
cbo	coupling between objects	increased when the methods of one class access services of another.
ce	efferent couplings	how many other classes is used by the specific class.
dam	data access	ratio of the number of private (protected) attributes to the total number of attributes
dit	depth of inheritance tree	
ic	inheritance coupling	number of parent classes to which a given class is coupled (includes counts of methods and variables inherited)
lcom	lack of cohesion in methods	number of pairs of methods that do not share a reference to an instance variable.
lcom3	another lack of cohesion measure	if m, a are the number of <i>methods, attributes</i> in a class number and $\mu(a)$ is the number of methods accessing an attribute, then $lcom3 = ((\frac{1}{a} \sum_j^a \mu(a_j)) - m)/(1 - m)$.
loc	lines of code	
max_cc	maximum McCabe	maximum McCabe's cyclomatic complexity seen in class
mfa	functional abstraction	number of methods inherited by a class plus number of methods accessible by member methods of the class
moa	aggregation	count of the number of data declarations (class fields) whose types are user defined classes
noc	number of children	
npm	number of public methods	
rfc	response for a class	number of methods invoked in response to a message to the object.
wmc	weighted methods per class	
defects	defects	number of defects per class, seen in post-release bug-tracking systems.

Fig. 8. OO measures used in the LUCENE, XALAN, JEDIT, VELOCITY, SYNAPSE, TOMCAT, XERCES defect data sets. The last line shows the dependent attribute (defects reported to a post-release bug-tracking system).

that is half of the maximum value seen in the raw data D_{all} .

- Find the 25, 50, 75, 100th percentile normalized value in $D_{treated}$ and D_{all} . For both distributions, let the *median, stability,* and *worst case* values to be the 50th, 75th-25th, and 100th percentile values, respectively.

Formally, in each case, the following steps were done:

1. Combine all the data from all sources; e.g., for defect reduction, combine together data from $S_1 = \text{JEDIT}$, $S_2 = \text{XERCES}$, $S_3 = \text{XALAN}$, etc.
2. Cluster the combined data using WHERE. Note that each cluster may now contain examples from multiple sources.
3. For each cluster C , find C' (the neighboring cluster with the best median score of its dependent variable).
4. Apply WHICH to C' to learn some rules. then apply those learned rules, i.e., that predicts for lower defects, to C .
5. Report the *median, stability,* and *worst case* values in *all* and *some*.

For experiments 2.a and 2.b, step 3 becomes:

- 3ab. For data in C from source S_i , find examples in C' from other sources S_j , where $i \neq j$.

For experiment 2.a, step 4 becomes:

- 4a. Apply WHICH to the S_j data from C' . In this experiment, the S_i examples in C are the *all* set

while *some* are the examples selected by the rules learned from other sources S_j in C' .

For the experiment 2.b, step 4 becomes:

- 4b. The S_i examples in C are treated with the rules learned from the same sources S_i in C' .

Next, we discuss the details of the experiments separately along with their results.

4.1 Experiment #1: Global Rules are Sub-Optimal

In the first experiment we compared *global lessons* with *cluster lessons*. We used WHERE to cluster all the available data, separately for effort estimation and defect prediction. Our clustering algorithm takes care to skip the dependent variable during clustering. That is, when computing the distance between examples (i.e., projects for effort estimation and classes for defect prediction), the defect or effort values in those examples are ignored. The dependent variable is used only after clustering.

Once WHERE divides the data, WHICH is applied to each cluster C to find the *cluster* rules for selecting lower defect/effort values. Next, the dependent values are used to score each cluster (median defect/effort values of projects in that cluster). For each cluster C we search its neighbors for the cluster $C' \neq C \wedge \text{adjacent}(C, C') \wedge (\text{score}(C') < \text{score}(C))$ with the best score (lowest median defect/effort values). The rule $C'.cluster$ is then applied to C to find

cluster	effort		defect						
	NASACOC	CHINA	LUCENE	XALAN	JEDIT	VELOCITY	SYNAPSE	TOMCAT	XERCES
global	kloc=1	afp=1	rfc=2	loc=1	rfc=2	cam=7	amc=1	loc=2	cbo=1
C0									
C1	rely=n	added=4	amc=7	amc=1	ic=7	noc=1	dit=4	cbm=1	dit=1
C2	prec=h	deleted=1	ca=1	cam=2	noc=1	dam=1 or 5		dam=1	dam=1
C3		deleted=1	dam=5	cam=3	amc=6	avg_cc=4		noc=1	ca=1 or 7
C4			mfa=1	dit=2 or 4	noc=1	moa=1		rfc=5	<u>cbo=1</u>
C5			moa=1	<u>loc=1</u>				lcom3=5	
C6				<u>loc =1 or 2</u>				max_cc=1	
C7				moa=1				cbm=1	

Fig. 9. *Global* and *cluster* rules for each cluster C_i . In this study, all the numerics were discretized into seven equal frequency bins so, e.g., kloc=1 should be read as “set kloc to its minimal value”. The **underlined bold entries** denote the rules that were the same globally as at cluster level. Note that there are very few rules that are the same globally as at cluster level.

		effort		defect							$\frac{\sum_{global}}{\sum_{cluster}}$	
		NASACOC	CHINA	LUCENE	XALAN	JEDIT	SYNAPSE	VELOCITY	TOMCAT	XERCES		
median	global	17	4	3	12	0	0	0	0	0	0	0.64
(50th percentile)	cluster	7	3	0	12	0	0	0	0	0	1	
stability	global	16	7	10	12	0	11	8	0	1	1	0.37
(75th-25th percentile)	cluster	6	6	3	0	0	0	8	0	1	1	
worst-case	global	100	100	100	100	100	100	100	100	100	100	0.39
(100th percentile)	cluster	9	100	23	62	8	33	50	33	32	32	

Fig. 10. Global vs. cluster reasoning. In this figure, *smaller* values are *better*. All values are percentages of the maximum effort or defect values seen in the untreated data sets so, e.g., a median global value of “17” in NASACOC is 17% of the maximum effort value in the NASACOC data sets.

$treated \subseteq C$. See above the definition for *treated* and *score*.

For every cluster C we compared the two treated sets:

- One using $C'.cluster$ and
- The other using rules generated globally across all the data (by applying WHICH to all the data sets, rather than to a single cluster).

We found that the rules generated from cluster lessons were better and different than those learned globally. Figure 9 shows the rules learned in different data sets. Each data set generated between 2 (SYNAPSE) to 8 (XALAN) clusters. One cluster always had the best score (lowest effort or defects) and this cluster was labeled C_0 . No rules were learned for this cluster since our recommendation for projects in C_0 is to “maintain the status quo” - i.e., we do not know how to improve the median value of the dependent variable (defect or effort) for this cluster, given the current data. Lines $C_i \in C_1..C_7$ show the cluster rules learned from C_i 's best neighbor. The underlined cluster rules are those that are same as the global rules (these appear in the results for XALAN and XERCES). Note that there are very few cluster rule sets that are the same as the global rules.

The effects of applying these rules are shown in Figure 10. These results are expressed in terms of percentile bands: *median*, *stability*, and *worst-case*. All values are expressed as ratios of maximum values

seen in the untreated data set (e.g., “50” means the middle value of the untreated data).

The first thing to note in Figure 10 is that our rule learning method is effective: the median and stability values of the above are small percentages of original data. Indeed, in five of the seven defect prediction experiments, the cluster rules selected projects with zero defects.

The second important conclusion of Figure 10 is that the cluster rules performed better than the global rules:

- The *median* performance of cluster rules is significantly better than the median performance of the global rules (Wilcoxon, 95% confidence).
- The *stability* around the median is greater with cluster rules than with global rules, i.e., if we use the global rules, then we will be *less confident* in our predictions on the effects of those rules.
- The *worst-case* results with global rules are far worse than the worst-case results of cluster rules. In all cases, the worst-case performance of the global rules is the same as the untreated data (see the second last row of the table). On the other hand, in seven results, the worst-case performance of the cluster rules is one-third or less of the maximum in the untreated data.

Why is the performance of the global rules sub-optimal? Our hypothesis is that software construction is such a diverse task that any global conclusion that

			VELOCITY1.4	POI1.5, ANTI1.6	poi-3.0, VELOCITY1.6, XERCES1.3	ANTI1.5, XERCES1.4, VELOCITY1.4	POI2.5, JEDIT3.2, LUCENE2.4	ANTI1.4, IVY1.4, XERCES1.2	IVY1.4, ANTI1.4, JEDIT4.3	POI3.0, JEDIT4.0, ivy-1.4	POI2.5, XERCES1.3, ANTI1.5	
median	within	5	0	2	4	0	0	0	0	0	0	0.4
(50th percentile)	cross	0	0	2	2	0	0	0	0	0	0	
stability	within	4	0	2	0	0	0	0	0	0	0	1.3
(75th-25th percentile)	cross	0	0	0	0	0	0	0	4	8	0	
worst case	within	10	32	15	16	33	67	9	23	0.3		
(100th percentile)	cross	0	11	2	11	33	0	4	3			

Fig. 11. Experiment 2. Using groups of three projects, S_j (the test cluster) is picked at random (marked in **bold**).

holds *across all projects* may be somewhat different to the conclusions learned from *individual projects*.

In summary, the data we obtained is not supportive of the claim that global lessons are the best for defect prediction or effort estimation. However, it would be a mistake to conclude that this experiment supports delphi localisms. As shown in our next experiment, the best way to divide the data is not via delphi localizations (into, say, just the JEDIT or LUCENE data sets) but by automatic localizations that cross data set boundaries.

4.2 Experiment #2: Local Rules are Sub-Optimal

IF a data miner is overly zealous, then it could read too much into the training data. In this situation, the learner may *overfit* its models to spurious details in the training data. Hence, many data mining algorithms employ overfitting-avoidance strategies to prune away needless elaborations. For example:

- Early versions of decision tree learners produced very big and bushy decision trees. One of the key innovations of C4.5 [61] was *leaf pruning* which prunes leaves until the error rate in the pruned tree starts increasing.
- After the INDUCT rule learner builds a rule condition, it runs a greedy back-select algorithm that checks if any condition can be removed, while preserving the accuracy of that rule [62].

Accordingly, the following experiment checked if WHERE and WHICH can be improved by an overfitting-avoidance strategy. Specifically, we selected the data used for training in each cluster, based on which source they come from (as explained in the beginning of this section - see step 3ab). This resulted in two ways of learning the rules.

We call the steps selecting data from different sources the *cross* treatment since it ensures that data from one source are treated with rules learned from data from other sources. We call the generated rules *neighbor rules*, because the source of the data is not random, as it comes from a neighboring cluster. The second treatment is referred to as *within* treatment and the rules are called *local rules*, as the data are treated with the rules learned from the same sources (caveat: providing that the data falls into separate neighboring C, C' clusters). It is important to note that in both treatments, rules are applied to the same subset S_i of C and that these rules select only some subset of that data.

As before, our results are expressed in terms of *median*, *stability*, and *worst-case*. We express those percentiles as a ratio of the maximum value seen in the untreated S_i data of C (i.e., the maximum value seen before any rule selects some subset). That is to say, all our results are expressed as values ranging from 0% to 100%.

The results are shown in Figure 11. In order to ensure external validity, the above procedure is repeated multiple times. The results of each repetition are shown in different columns of Figure 11. Each repetition combines together the data from three sources. Within each group of three, one source is the designated test cluster S_j (again, selected at random). While the data sets were picked at random from the PROMISE data repository, all had to conform to the same ontology (in practice, that meant that Experiment #2 was conducted on data sets of the OO ontology of Figure 8, since this is the most frequently shared ontology in PROMISE).

One issue that arises in these results is that the *within* treatment is so effective that it is sometimes

hard to distinguish further improvements. For example, at first glance, the *stability* results of *cross* seem worse than *within*. However, the raw values for *stability* are so small that a single value (the “8” in the last column) can throw off the results. The same effect (that the baseline *within* results are very small) also confuses an analysis of the *median* results.

However, when we turn to the *worst case* results, the numbers are large enough to allow for a differentiation of the *within* and *cross* results. We observe that the *worst case* results of the *cross* study are much better than the *within*.

In summary, based on an analysis of rules learned from neighboring clusters in different sources, we conclude that it is sub-optimal to learn purely local rules from the clusters within the same source as the test data. This is the standard overfitting-avoidance result [54]: it is best not to learn too much from local data. Rather, it is better for a learner to step back a little and train from related concepts (rather than concepts that are too similar).

We conjecture that the diverse nature of software construction means that even for projects built within the same organization, it is more useful to chase external data sources than just to use the local historical data. However, when using data from other sources, it is best to cluster that data and just use a small portion of data from the nearest cluster.

5 EXTERNAL VALIDITY

EXTERNAL validity is the issue of the generality of the conclusions of a study to data not used in that particular study. Within the specific context of effort estimation and defect prediction, the central claim of this paper is that any standard discussion on external validity cannot be trusted for SE data. That is, the methodology of this paper precludes a claim of external validity about specific conclusions (e.g., the relationship of inheritance depth to defects).

But this paper is not a counsel of despair. An essential feature of our work is that the same algorithms were used to generate recommendations for *both* defect reduction *and* effort reduction. This makes this paper somewhat unique since, in the literature, effort estimation and defect prediction are usually explored by different research teams² and techniques. That is:

- While this paper has shown that any *specific conclusion* about reducing effort or defects may not be externally valid (since they are project dependent) ...
- ... our externally valid *meta-conclusion* is that there are general techniques (i.e., WHERE+WHICH) for finding local conclusions in different projects.

2. Exception- see the work of Martin Shepperd who explores both areas [63], [64].

In other words, while we do not provide general lessons that work for defect prediction or effort estimation on any project, we provide instead a technique (using WHERE and WHICH) that can be used on any project with the available data.

To disprove the external validity of this meta-conclusion, a research team would need to demonstrate the stability of conclusions across multiple projects. We would propose one sanity check for that demonstration: the project-independent conclusion must be “significant”.

For examples of less-than-significant conclusions, we refer the reader to the *global* results shown in the first line of Figure 9. Here, we read that effort or defects can be reduced by:

- minimizing lines of code: see the *kloc=1* results for NASACOC and *loc = 1* for XALAN;
- minimizing function points: see the *afp=1* result from CHINA.

That is, those global conclusions are just trite statements that if programmers write less code, they can do so in less time (and introduce fewer defects). In most development projects, programmers do not have the option of writing less code. What makes WHICH’s local conclusions non-trite is that most often they are not about merely reducing the size of a system. Rather, as seen in Figure 9, they are about system reorganization (such as reorganizing the class hierarchy).

Note that, in terms of the current literature, current results support the external validity of the conclusions of this paper. Inspired by our earlier work [10], Bettenburg et al. [6] recently repeated our experiment #1 using different clustering tools. That study found the same conclusion as ours, i.e., that cluster rules do better than rules learned across the whole data.

6 CONCLUSION AND FUTURE WORK

THIS paper has discussed ways to collect training data for learning lessons from SE projects and to generate rules for reducing the number of defects or the development effort.

At issue here is how the chief information officers should set policies for their projects. Should they demand that all their projects rigorously apply all the advice from the standard SE textbooks or other literature? Or should they devote resources to a “local lessons team” that explores local data to find the best local practices?

The results of this paper strongly endorse the creation of the local lessons team. However, that team should apply automatic algorithms to build clusters from all available data. The best cluster is the one that is nearby (see Experiment #1) but not from the same source as the test data (see Experiment #2). While we experimented only on defect prediction and effort estimation data, our procedure is data agnostic and

we believe it would apply to other kinds of data, for other tasks.

More generally, in terms of SE research, the experiments of this paper show that a software engineering project is an intricate entity that is best described in terms of a complex combination of multi-dimensional factors. Hence:

- Trite global rules are not sufficient for controlling such complex entities, at least when it comes to defect prediction and effort estimation.
- It is neither sufficient to characterize the data with simple divisions of data into local contexts. Before researchers attempt to draw lessons from some historical data source, they should (a) ignore any existing local divisions into multiple sources; (b) cluster across all available data; then (c) restrict the learning of lessons to the clusters from other sources that are nearest to the test data.

As to future work, it is important to check how often, in other data sets, cluster rules are better than global rules.

Other future work might include improving WHICH and WHERE. WHICH contains numerous design options that deserve closer attention. Like any beam search, WHICH only stores the top ($N = 50$) rules in its stack. Also, prior to building rules, WHICH discretizes numeric data into $B = 7$ equal frequency bins. It is possible that different values of N and B would result in better rules. More fundamentally, WHICH is one of a large class of contrast set learners. Other candidate contrast set learners, which might do better than WHICH, are MINWAL [65], STUCCO [66], and others [67].

Another task deserving future research is to explore multi-goal optimization. In the results of Figure 10, we tried to minimize effort in the NASACOC and CHINA data sets while minimizing defects in the remaining seven data sets. A more challenging goal would be to reduce defects and effort at the same time. The search-based software engineering literature [68] lists many techniques that might be useful in this regard, including tabu search [69], ant algorithms [70], and particle swarm optimization [71].

7 ACKNOWLEDGEMENTS

The work was funded by NSF grants CCF:1017330 and CCF: 1017263, the Qatar/WVU research grant NPRP 09-12-5-2-470, partly funded by the Finnish Funding Agency for Technology and Innovation (TEKES) under Cloud Software Program and the Academy of Finland with Grant Decision No. 260871. Also, the research reported in this document was performed in connection with contract/instrument W911QX-10-C-0066 with the U.S. Army Research Laboratory. The views and conclusions contained in this document/presentation are those of the authors and

should not be interpreted as presenting the official policies or position, either expressed or implied, of the U.S. Army Research Laboratory or the U.S. Government unless so designated by other authorized documents. Citation of manufacturers or trade names does not constitute an official endorsement or approval of the use thereof.

REFERENCES

- [1] C. Jones, *Estimating Software Costs*, 2nd Edition. McGraw-Hill, 2007.
- [2] B. Boehm, E. Horowitz, R. Madachy, D. Reifer, B. K. Clark, B. Steece, A. W. Brown, S. Chulani, and C. Abts, *Software Cost Estimation with Cocomo II*. Prentice Hall, 2000.
- [3] B. A. Kitchenham, T. Dyba, and M. Jorgensen, "Evidence-based software engineering," in *ICSE '04: Proceedings of the 26th International Conference on Software Engineering*. Washington, DC, USA: IEEE Computer Society, 2004, pp. 273–281.
- [4] N. Nagappan, T. Ball, and A. Zeller, "Mining metrics to predict component failures," in *Proceedings of the 28th international conference on Software engineering*, ser. ICSE '06, 2006, pp. 452–461.
- [5] D. Posnett, V. Filkov, and P. Devanbu, "Ecological inference in empirical software engineering," in *Proceedings of ASE'11*, 2011.
- [6] N. Bettenburg, M. Nagappan, and A. E. Hassan, "Think locally, act globally: Improving defect and effort prediction models," in *MSR'12*, 2012, pp. 60–69.
- [7] Y. Yang, L. Xie, Z. He, Q. Li, V. Nguyen, B. W. Boehm, and R. Valerdi, "Local bias and its impacts on the performance of parametric estimation models," in *PROMISE'11*, 2011.
- [8] K. Lum, J. Hihn, and T. Menzies, "Studies in software cost model behavior: Do we really understand cost model performance?" in *ISPA Conference Proceedings*, 2006, available from <http://menzies.us/pdf/06ispa.pdf>.
- [9] T. Menzies, J. Greenwald, and A. Frank, "Data mining static code attributes to learn defect predictors," *IEEE Transactions on Software Engineering*, January 2007, available from <http://menzies.us/pdf/06learnPredict.pdf>.
- [10] T. Menzies, A. Butcher, A. Marcus, T. Zimmermann, and D. Cok, "Local vs global models for effort estimation and defect prediction," in *IEEE ASE'11*, 2011, available from <http://menzies.us/pdf/11ase.pdf>.
- [11] T. Menzies, Z. Chen, D. Port, and J. Hihn, "Simple software cost estimation: Safe or unsafe?" in *Proceedings, PROMISE workshop, ICSE 2005*, 2005, available from <http://menzies.us/pdf/05safewhen.pdf>.
- [12] B. Boehm, *Software Engineering Economics*. Prentice Hall, 1981.
- [13] H. Olague, L. Etzkorn, S. Gholston, and S. Quattlebaum, "Empirical validation of three software metrics suites to predict fault-proneness of object-oriented classes developed using highly iterative or agile software development processes," *Software Engineering, IEEE Transactions*, vol. 33, no. 6, pp. 402–419, 2007.
- [14] K. Aggmarwal, Y. Singh, A. Kaur, and R. Malhotra, "Empirical analysis for investigating the effect of object-oriented metrics on fault proneness: a replicated case study," *Software Process: Improvement and Practice*, vol. 14, no. 1, January 2009.
- [15] E. Arisholm and L. Briand, "Predicting fault prone components in a JAVA legacy system," *2006 ACM/IEEE international symposium on Empirical software engineering*, p. 17, 2006.
- [16] V. Basili, L. Briand, and W. Melo, "A validation of object-oriented design metrics as quality indicators," *Software Engineering, IEEE Transactions*, vol. 22, no. 10, pp. 751–761, 1996.
- [17] L. Briand, J. Wust, J. Daly, and D. V. Porter, "Exploring the relationships between design measures and software quality in object-oriented systems," *Journal of Systems and Software*, vol. 51, no. 3, pp. 245–273, 2000.
- [18] L. Briand, J. Wust, and H. Lounis, "Replicated case studies for investigating quality factors in object-oriented designs," *Empirical Software Engineering*, vol. 6, no. 1, pp. 11–58, 2001.

- [19] M. Cartwright and M. Shepperd, "An empirical investigation of an object-oriented software system," *Software Engineering, IEEE Transactions*, vol. 26, no. 8, pp. 786–796, 2000.
- [20] K. el Emam, W. Melo, and J. Machado, "The prediction of faulty classes using object-oriented design metrics," *Journal of Systems and Software*, vol. 56, no. 1, pp. 63–75, 2001.
- [21] K. el Emam, S. Benlarbi, N. Goel, and S. Rai, "A validation of object-oriented metrics," *National Research Council of Canada, NRC/ERB*, vol. 1063, 1999.
- [22] M. Tang, M. Kao, and M. Chen, "An empirical study on object-oriented metrics," *Software Metrics Symposium*, vol. Proceedings. Sixth International, pp. 242–249, 1999.
- [23] P. Yu, T. Systa, and H. Muller, "Predicting fault-proneness using oo metrics an industrial case study," *Sixth European Conference on Software Maintenance and Reengineering*, pp. 99–107, 2002.
- [24] R. Subramanyam and M. S. Krishnan, "Empirical analysis of ck metrics for object-oriented design complexity: implications for software defects," *IEEE Transactions on Software Engineering*, vol. 29, no. 4, pp. 297–310, April 2003.
- [25] Y. Zhou and H. Leung, "Empirical analysis of object-oriented design metrics for predicting high and low severity faults," *Software Engineering, IEEE Transactions*, vol. 32, no. 10, pp. 771–789, 2006.
- [26] T. Gyimothy, R. Ferenc, and I. Siket, "Empirical validation of object-oriented metrics on open source software for fault prediction," *Software Engineering, IEEE Transactions*, vol. 31, no. 10, pp. 897–910, 2005.
- [27] T. Holschuh, M. Pauser, K. Herzig, T. Zimmermann, R. Premraj, and A. Zeller, "Predicting defects in SAP Java code: An experience report," *Software Engineering - Companion Volume, 2009. ICSE-Companion 2009. 31st International Conference*, pp. 172–181, 2009.
- [28] R. Shatnawi and W. Li, "The effectiveness of software metrics in identifying error-prone classes in post-release software evolution process," *Journal of Systems and Software*, vol. 81, no. 11, pp. 1868–1882, 2008.
- [29] F. Fioravanti and P. Nesi, "A study on fault-proneness detection of object-oriented systems," *Software Maintenance and Reengineering, 2001. Fifth European Conference*, pp. 121–130, 2001.
- [30] M. Thongmak and P. Muenchaisri, "Predicting faulty classes using design metrics with discriminant analysis," *Software Engineering Research and Practice*, pp. 621–627, 2003.
- [31] G. Denaro, L. Lavazza, and M. Pezze, "An empirical evaluation of object oriented metrics in industrial setting," *The 5th CaberNet Plenary Workshop, Porto Santo, Madeira Archipelago, Portugal*, 2003.
- [32] A. Janes, M. Scotto, W. Pedrycz, B. Russo, M. Stefanovic, and G. Succi, "Identification of defect-prone classes in telecommunication software systems using design metrics," *Information Sciences*, vol. 176, no. 24, pp. 3711–3734, 2006.
- [33] M. English, C. Exton, I. Rigon, and B. Cleary, "Fault detection and prediction in an open-source software project," *Proceedings of the 5th International Conference on Predictor Models in Software Engineering*, no. 1-11, 2009.
- [34] R. Shatnawi, "A quantitative investigation of the acceptable risk levels of object-oriented metrics in open-source systems," *IEEE Transactions on Software Engineering*, vol. 36, no. 2, pp. 216–225, 2010.
- [35] Y. Singh, A. Kaur, and R. Malhotra, "Empirical validation of object-oriented metrics for predicting fault proneness models," *Software Quality Journal*, vol. 18, no. 1, pp. 3–35, 2010.
- [36] D. Glasberg, K. el Emam, W. Memo, and N. Madhavji, "Validating object-oriented design metrics on a commercial JAVA application," *NRC 44146*, 2000.
- [37] K. el Emam, S. Benlarbi, N. Goel, and S. Rai, "The confounding effect of class size on the validity of object-oriented metrics," *Software Engineering, IEEE Transactions*, vol. 27, no. 7, pp. 630–650, 2001.
- [38] M. Thapaliyal and G. Verma, "Software defects and object oriented metrics-an empirical analysis," *International Journal of Computer Applications*, vol. 9/5, 2010.
- [39] J. Xu, D. Ho, and L. Capretz, "An empirical validation of object-oriented design metrics for fault prediction," *Journal of Computer Science*, pp. 571–577, July 2008.
- [40] G. Succi, "Practical assessment of the models for identification of defect-prone classes in object-oriented commercial systems using design metrics," *Journal of Systems and Software*, vol. 65, no. 1, pp. 1–12, January 2003.
- [41] M. Jorgensen, "A review of studies on expert estimation of software development effort," *Journal of Systems and Software*, vol. 70, no. 1-2, pp. 37–60, 2004.
- [42] 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, member-Kitchenham, Barbara A.
- [43] S. G. MacDonell and M. J. Shepperd, "Comparing local and global software effort estimation models – reflections on a systematic review," in *Proceedings of the First International Symposium on Empirical Software Engineering and Measurement*, ser. ESEM '07. Washington, DC, USA: IEEE Computer Society, 2007, pp. 401–409.
- [44] C. Mair and M. Shepperd, "The consistency of empirical comparisons of regression and analogy-based software project cost prediction," in *Empirical Software Engineering, 2005. 2005 International Symposium on*, 2005, p. 10 pp.
- [45] T. Zimmermann, N. Nagappan, H. Gall, E. Giger, and B. Murphy, "Cross-project defect prediction," in *ESEC/FSE'09*, August 2009.
- [46] L. C. Briand and J. Wuest, "Empirical studies of quality models in object-oriented systems," *Advances in Computers*, vol. 56, no. 0, 2002, edited by M. Zelkowitz.
- [47] B. Turhan, "On the dataset shift problem in software engineering prediction models," *Empirical Software Engineering*, vol. 17, pp. 62–74, 2012.
- [48] T. Menzies, Z. Milton, B. Turhan, B. Cukic, Y. Jiang, and A. Bener, "Defect prediction from static code features: Current results, limitations, new approaches," *Automated Software Engineering*, no. 4, December 2010, available from <http://menzies.us/pdf/10which.pdf>.
- [49] G. Gay, T. Menzies, M. Davies, and K. Gundy-Burlet, "Automatically finding the control variables for complex system behavior," *Automated Software Engineering*, no. 4, December 2010, available from <http://menzies.us/pdf/10tar34.pdf>.
- [50] T. Menzies and Y. Hu, "Data mining for very busy people," in *IEEE Computer*, November 2003, available from <http://menzies.us/pdf/03tar2.pdf>.
- [51] 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>.
- [52] Q. Du and J. E. Fowler, "Low-Complexity Principal Component Analysis for Hyperspectral Image Compression," *International Journal of High Performance Computing Applications*, vol. 22, no. 4, pp. 438–448, Nov. 2008.
- [53] C. Faloutsos and K.-I. Lin, "Fastmap: a fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets," in *Proceedings of the 1995 ACM SIGMOD international conference on Management of data*, ser. SIGMOD '95. New York, NY, USA: ACM, 1995, pp. 163–174. [Online]. Available: <http://doi.acm.org/10.1145/223784.223812>
- [54] I. H. Witten and E. Frank, *Data mining. 2nd edition*. Los Altos, US: Morgan Kaufmann, 2005.
- [55] D. Rumelhart, G. Hinton, and R. Williams, "Learning representations by back-propagating errors," *Nature*, no. 323, pp. 533 – 536, 1986.
- [56] S. Bouktif, H. Sahraoui, and G. Antoniol, "Simulated annealing for improving software quality prediction," in *GECCO '06: Proceedings of the 8th annual conference on Genetic and evolutionary computation*, 2006, pp. 1893–1900.
- [57] J. R. Quinlan, "Learning with Continuous Classes," in *5th Australian Joint Conference on Artificial Intelligence*, 1992, pp. 343–348, available from <http://citeseer.nj.nec.com/quinlan92learning.html>.
- [58] G. Boetticher, "When will it be done? the 300 billion dollar question, machine learner answers," *IEEE Intelligent Systems*, June 2003.
- [59] E. Kocaguneli, G. Gay, T. Menzies, Y. Yang, and J. W. Keung, "When to use data from other projects for effort estimation," in *IEEE ASE'10*, 2010, available from <http://menzies.us/pdf/10other.pdf>.

- [60] T. Menzies, B. Caglayan, E. Kocaguneli, J. Krall, F. Peters, and B. Turhan. (2012, June) The promise repository of empirical software engineering data. [Online]. Available: <http://promisedata.googlecode.com>
- [61] R. Quinlan, *C4.5: Programs for Machine Learning*. Morgan Kaufman, 1992, ISBN: 1558602380.
- [62] B. Gaines and P. Compton, "Induction of ripple down rules," in *Proceedings, Australian AI '92*. World Scientific, 1992, pp. 349–354.
- [63] Q. Song, Z. Jia, M. Shepperd, S. Ying, and J. Liu, "A general software defect-proneness prediction framework," *IEEE Trans. Softw. Eng.*, vol. 37, pp. 356–370, May 2011. [Online]. Available: <http://dx.doi.org/10.1109/TSE.2010.90>
- [64] M. Jorgensen and M. Shepperd, "A systematic review of software development cost estimation studies," *IEEE Trans. Softw. Eng.*, vol. 33, pp. 33–53, January 2007. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1248721.1248736>
- [65] C. Cai, A. Fu, C. Cheng, and W. Kwong, "Mining association rules with weighted items," in *Proceedings of International Database Engineering and Applications Symposium (IDEAS 98)*, August 1998, available from http://www.cse.cuhk.edu.hk/~kdd/assoc_rule/paper.pdf.
- [66] S. Bay and M. Pazzani, "Detecting change in categorical data: Mining contrast sets," in *Proceedings of the Fifth International Conference on Knowledge Discovery and Data Mining*, 1999, available from <http://www.ics.uci.edu/~pazzani/Publications/stucco.pdf>.
- [67] P. K. Novak, N. Lavrač, and G. I. Webb, "Supervised descriptive rule discovery: A unifying survey of contrast set, emerging pattern and subgroup mining," *J. Mach. Learn. Res.*, vol. 10, pp. 377–403, June 2009.
- [68] M. Harman and B. Jones, "Search-based software engineering," *Journal of Information and Software Technology*, vol. 43, pp. 833–839, December 2001.
- [69] F. Glover and M. Laguna, "Tabu search," in *Modern Heuristic Techniques for Combinatorial Problems*, C. Reeves, Ed. Oxford, England: Blackwell Scientific Publishing, 1993. [Online]. Available: citeseer.ist.psu.edu/glover97tabu.html
- [70] M. Zlochin, M. Birattari, N. Meuleau, and M. Dorigo, "Model-based search for combinatorial optimization: A critical survey," *Annals of Operations Research*, vol. 131, pp. 373–395, 2004.
- [71] H. Pan, M. Zheng, and X. Han, "Particle swarm-simulated annealing fusion algorithm and its application in function optimization," in *International Conference on Computer Science and Software Engineering*, 2008, pp. 78–81.

Tim Menzies (Ph.D., UNSW) is a Professor in CSEE at WVU and the author of over 200 referred publications. At WVU, he has been a lead researcher on projects for NSF, NIJ, DoD, NASA's Office of Safety and Mission Assurance, as well as SBIRs and STTRs with private companies. He teaches data mining and artificial intelligence and programming languages. Tim is the co-founder of the PROMISE conference series devoted to reproducible experiments in software engineering. He is an associate editor of IEEE Transactions on Software Engineering and the Automated Software Engineering Journal. In 2012, he served as co-chair of the program committee for the IEEE Automated Software Engineering conference.

Andrew Butcher is an active member of the open source community and a software engineering at Redhat (East Carolina). For further information, see his Git repositories <https://github.com/abutcher> or his home page <http://abutcher.org/>.

David Cok is Associate Vice President of Technology at GrammaTech. He graduated with honors from Calvin College with A.B. degrees in Physics and Mathematics. He also earned an A.M. and Ph.D. in Physics from Harvard University. Dr. Cok's career includes research in static analysis with an emphasis on practical application to industrial software development. Before joining GrammaTech, Dr. Cok was a Senior Research Associate and leader of research and product development groups at Eastman Kodak Company.

Andrian Marcus is an Associate Professor in the Department of Computer Science at Wayne State University (Detroit, MI), where he is faculty since 2003. He obtained his Ph.D. in Computer Science from Kent State University and he has prior degrees from The University of Memphis (Memphis, TN) and Babes-Bolyai University (Cluj-Napoca, Romania). His current research interests are in software engineering, with focus on software evolution and program comprehension. He is best known for his work on using information retrieval and text mining techniques for software analysis to support comprehension during software evolution. He is a member of the Steering Committee of the IEEE International Conference on Software Maintenance and he served as the conference's General Chair and Program Co-Chair.

Lucas Layman Dr. Layman's past research includes reliability and productivity prediction at Microsoft Corporation, safety measurement at NASA, and agile process engineering and evaluation at the Department of Defense, IBM, and Sabre Airline Solutions. In addition to empirical studies on process and process improvement, Dr. Layman has investigated human factors in software development and authored a number of papers on computer science education. He received his Ph.D. in Computer Science from North Carolina State University in 2009 and his B.S. in Computer Science from Loyola College in 2002. In his current position, Lucas is a research scientist at the Fraunhofer center for experimental software engineering, Maryland, USA.

Forrest Shull is a senior scientist at the Fraunhofer Center for Experimental Software Engineering in Maryland (CESE), a non-profit research and tech transfer organization, where he leads the Measurement and Knowledge Management Division. At Fraunhofer CESE, he has been a lead researcher on projects for NASA's Office of Safety and Mission Assurance, the NASA Safety Center, the U.S. Department of Defense, the National Science Foundation, the Defense Advanced Research Projects Agency (DARPA), and companies such as Motorola and Fujitsu Labs of America. He is an associate adjunct professor at the University of Maryland College Park, and Editor-in-Chief of IEEE Software.

Burak Turhan is a Tutkijatohtori (Doctor researcher) at the Department of Information Processing Science, University of Oulu in Finland, and a member of the PROMISE steering committee. For further information, see his website: <http://turhanb.net>

Thomas Zimmermann received his Ph.D. degree from Saarland University in Germany. He is a researcher in the Research in Software Engineering group at Microsoft Research Redmond, adjunct assistant professor at the University of Calgary, and affiliate faculty at University of Washington. His research interests include empirical software engineering, mining software repositories, computer games, recommender systems, development tools, and social networking. He is best known for his research on systematic mining of version archives and bug databases to conduct empirical studies and to build tools to support developers and managers.