

¹Venkata Sai Rahul
Trivedi Kothapalli

XTREE: An AI-Driven Framework for Planning Effective Software System Improvements



Abstract: - Business decision-makers increasingly demand actionable insights from analytics, moving beyond predictive analysis to decision-making support. This paper presents XTREE, an artificial intelligence (AI)-based decision support framework for software projects. XTREE employs supervised machine learning, specifically decision tree algorithms, to learn and recommend plans that improve software metrics such as defect counts and runtime efficiency. By analyzing datasets with weighted class labels indicating quality, XTREE suggests precise changes to input features to transition a software system from a “bad” to a “better” quality state. Evaluated on 11 datasets across various software projects, XTREE outperformed three state-of-the-art planning methods, achieving median improvements of up to 56% and maximum improvements of 77% in defect reduction and runtime optimization, respectively. This research highlights XTREE’s effectiveness in delivering interpretable and actionable plans, making it a valuable tool for AI applications in software engineering and business analytics.

Keywords: XTREE, AI, Prediction, Planning

I. INTRODUCTION

Business users are demanding tools that support business-level interpretations of their data. At a panel on software analytics at ICSE’12, industrial practitioners lamented the state of the art in software analytics [1]. Panelists commented “prediction is all well and good, but what about decision making?”. Note that these panelists were more interested in the interpretations and follow-up that occurs after the mining, rather than just the mining itself. So:

- Instead of just accepting predictions on how many software defects to expect, business users might now demand a plan to reduce the likelihood of those defects.
- Instead of just accepting predictions on the runtime time of their software, business users might now demand a plan to reduce that runtime.

In response to this business-level demands for planners, we propose a novel planning method called XTREE for learning changes to a software system such that its performance “improves”, according to some measure. This paper uses XTREE to reduce the expected value of the defects in Jureczko et al.’s JAVA systems [2]; and the runtimes in software configured by Siegmund et al. [3]. The contributions of this paper are (1) the new XTREE algorithm and (2) an evaluation strategy that shows XTREE performing significantly “better” than planners proposed in our prior work [4], [5], where “better” means effective (plans change the expected values of the class); succinct (it is inconvenient if plans always require changing everything); and surprising (a planner should sometimes tell us things we do not expect since, otherwise, there is no value added from using the planner). The rest of this paper describes our data, our planners, and the experiment that ranks XTREE against alternate approaches. This is followed by notes on related work and validity.

¹Colorado State University Systems Engineering Department, Fort Collins, CO, USA
3waythe@gmail.com

II. PRELIMINARIES

- A. *What is a “Plan”?*
- B. Our planners use tables of data with independent features and a dependent class feature. Classes have weights that indicate what rows are “bad” or “better”. Plans change a row such that it is more likely to be “better”. Specifically, for every test example Z , planners proposes a plan Δ to adjust feature Z_j :

$$\forall \delta_j \in \Delta : Z_j = \begin{cases} Z_j + \delta_j & \text{if } Z_j \text{ is numeric} \\ \delta_j & \text{otherwise} \end{cases}$$

For example, to simplify a large bug-prone method, our planners might suggest to a developer to reduce its size (i.e. refactor that code by, say, splitting it across two simpler functions).

Note that we make no assumption that a plan mentions every feature (so plan1 can be more succinct than plan2 when plan1 mentions fewer features than plan2).

- C. *From Prediction to Planning*

This paper is about the next step after prediction. Suppose a business user is presented with a prediction and they do not like what they see; e.g. the runtimes are too long or the number of defects is too high. This user may then ask a planning question; i.e. “what can we change to do better than that?”.

Before exploring automatic methods to answer the planning question, we first comment on two manual methods.

One way to propose changes to a project would be to ask some smart experienced person for their opinion on how to (e.g.) reduce defects and/or decrease runtimes. Sometimes such advice is an effective strategy and sometimes it is not. According to Passos et al. [6], developers may assume that the lessons they learn from a few past projects are general to all their future projects. They comment “past experiences were taken into account without much consideration for their con- text” [6]. Jørgensen & Gruschke [7] offer a similar warning. They report that the supposed software engineering “Gurus” rarely use lessons from past projects to improve their future reasoning and that such poor past advice can be detrimental to new projects. [7]. Accordingly, we propose a “trust, but verify” approach. After a software guru offers some sage wisdom, it is wise to ask some other oracle if there are any better options (just as a sanity check). The rest of this paper discusses some methods to build automatic oracles to implement that sanity check.

Another way to find changes to a project might be to rely on the peer review processes used by the SE research community. This approach would propose changes to software projects that concur with internationally accept best practices. There are two problems with that approach. Firstly, given the rapid pace of change in software engineering, we may be asking questions for which there is no current widely accepted concept of “best practice”.

Secondly, given the diversity of SE products and practices and personnel, it may well be that the current project being discussed is substantively different to prior work. Numerous recent local learning results compare (1) single models learned from all available data to (2) multiple models learned from clusters within the data [8-14]. A repeated result in those studies is that the local models generated better effort and defect predictions (better median results, lower variance in the predictions). §IV-B3 of this paper offers yet another locality result:

- One standard rule in the literature is that it is useful to implement modules such that they are internally cohesive (use much of their own local methods) while being loosely coupled with other classes [15].
- While that may be true in general, for particular classes other changes may be more important (later in this paper, we show one set of results where that is indeed the case).

In summary, it is useful to have automatic methods to recom- mend changes. Such methods can fill in for human expertise (if such experts are absent) or to offer a second opinion. Also, prior to making automatic recommendations, it is wise to first stratify the data (clump it into related examples) then generate advice specific to each clump. Accordingly, the rest of this paper defines and evaluates automatic methods to find plans from N examples divided into many clumps.

D. *Trusting the Changes*

XTREE is evaluated by comparing predicted performance scores before and after a planner makes changes to the feature values of an example: After making those changes, we may have a new example that has never been seen before. Therefore, it must be asked “can we trust the predictions made on such new examples?”

To answer this question, we note that data miners explore two “clouds” of data: (1) the cloud of training examples and

(2) the cloud of test examples (for a visualization of these clouds, see Figure 2). We should mistrust the predictions made by a model if it is being applied to examples that are too far away from the training cloud. To test for “too far”, we can run a data mining experiment that tests how well a model learned from the training data applies to the test data. Such experiments return some performance value.

Note that predictions about changes that fall within the space of the training+test data, will be at least as accurate as the predictions of the original test data. With this, we assert that the predictions for changes that move examples towards/away from the training data can be trusted more/less (respectively).

Accordingly, we need trust-increasing planners to generate new examples closer to the training examples. To see how this works, Figure 2 is from the ivy data set (one of the Jureczko data sets used in this paper). It shows: (1) the training examples in gray, (2) the test examples in red, and (3) the changed examples displaced after applying a plan (in green). Note that the the changed examples cases (shown in green) fall closer to the training cases (shown in gray) than the test cases (shown in red).

In that green region of changed examples, our belief in the value of predictions will be just as much as, if not more than, our belief in the value of the predictions in the red region (that contains the original test data). This pattern of Figure 2 (where the new examples are found closer to the training cases than the test cases) has been observed in all the other data sets studied in this paper. Hence, we can assert that predictors learned from these training examples have some authority in the regions containing the changes examples.

That said, the above comes with some important caveats:

- The quality of the prediction depends on the nature of the training data. Thus, we strongly recommend that both the data set and the predictor be assessed prior to planning. This ensures that the predictor’s performance is adequate for a data set. We tackle this issue in detail in §??.
- Planners should be designed to be trust increasing. We list four such planning methods in §III.
- Where possible, planners should be assessed via some external oracle that can accurately assess new examples. For an example of that kind of analysis, see §V-C.

The Jureczko data records the number of known defects for each class using a post-release bug tracking system. The classes are described in terms of nearly two dozen metrics such as number of children (noc), lines of code (loc), etc. For details on the Jureczko data, see Figure 1 and Figure 3. All the planning methods of this paper reflect on the numeric value of the raw defect counts. The predictor considers defects to be a Boolean class data set, where defects are TRUE if the numeric defect count is greater than zero and FALSE otherwise.

The Siegmund data, described in Figure 4, records the

runtimes of compiled systems. To obtain the data, Siegmund et al. perturbed the configuration parameters in the Makefiles of six systems: Apache, SQLite, LLVM, x264 and two versions of the Berkeley database (one written in “C” and one in Java). Then, the performance was measured using standard bench- marking tools (delivered by ORACLE for Berkley data sets and other popular tools such as AUTOBENCH and HTTPREF for the rest) [3]. It’s worth noting that in each of the above data sets, several features are interdependent, this is expressed using a feature model. Figure 5 shows an example of such a feature model defining valid combinations of settings in the Berkley

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.
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	fferent 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 case variable.
lcom3	another lack of cohesion measure	if m, a are the number of methods and <i>attributes</i> in a class number and $\mu(a)$ is the number of methods accessing an attribute, then $lcom3 = ((1, j) \mu(a, j) - m) / (1 - m)$
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	
nDefects	raw defect counts	Numeric: number of defects found in post-release bug-tracking systems.
defect	defects present?	Boolean: if $nDefects > 0$ then <i>true</i> else <i>false</i>

Fig. 1. OO measures used in our defect data sets. Last lines, shown in denote the dependent variable

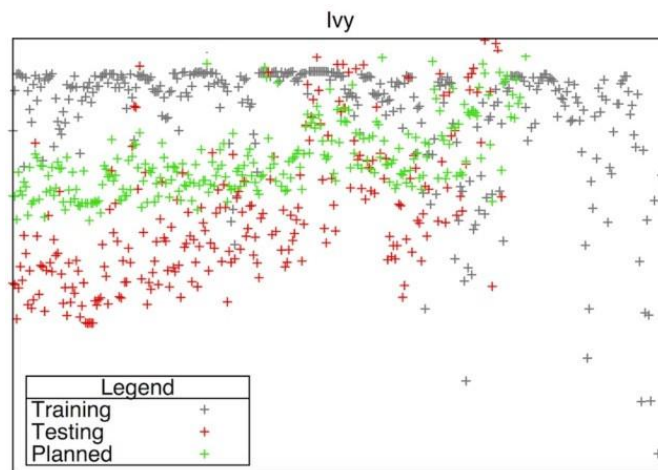


Fig. 2. Gray, red, green show (1) training examples, (2) test examples and (3) tests that have been altered by planners. This figure uses axes generated from the first two components of a PCA analysis of all points.

data set	cases	% defective
ant	947	22
camel	1819	19
jedit	1257	2
ivy	352	11
log4j	244	92
lucebe	442	59
poi	936	64
synapse	379	34
velocity	410	34
xalan	2411	99
xerces	1055	74

Fig. 3. Jureczko data: columns in the format of Figure 1.

These feature models were used by Siegmund et al. to ensure all their perturbations are valid (we will use the same models to cull invalid plans). Our evaluation strategy (discussed below) divides this data into a training and a testing set. From the training set we apply a data miner (to learn a quality predictor) and various planning methods (to learn different plans). Next, we try applying each of those plans to the test set and ask the quality predictor to assess the changed examples. Finally, we say that the “best” planner is the one that most reduces the predicted values in the changed examples.

Project	Domain	Lang.	LOC	Features	Config
BDBC: Berkeley DB	Database	C	219,811	18	2560
BDBJ: Berkeley DB	Database	Java	42,596	32	400
Apache	Web Server	C	230,277	9	192
SQLite	Database	C	312,625	39	3,932,160
LLVM	Compiler	C++	47,549	11	1024
x264	Video Enc.	C	45,743	16	1152

Fig. 4. Siegmund data



Fig. 5. Berkeley database feature model (“C” version).

data set	Data set properties					Results from learning								
	training		testing			tp	untuned		tuned		change			
	versions	cases	versions	cases	% defec		pf	good?	pd	pf	good?	pd	pf	
jedit	3.2, 4.0, 4.1, 4.2	1257	4.3	492		55	29	64	29	y	9	0	*	
ivy	1.1, 1.4	352	2.0	352		65	35	y	65	28	y	0	-7	*
camel	1.0, 1.2, 1.4	1819	1.6	965		4199	31		56	37		5	6	
ant	1.3, 1.4, 1.5, 1.6	947	1.7	745		49	13	y	63	16	y	14	3	*
synapse	1.0, 1.1	379	1.2	256		4354	19		47	15		2	-4	
velocity	1.4, 1.5	410	1.6	229		7384	60		76	60		-2	0	
lucene														

poi	2.0, 2.2	442	2.4	340	5569	25	60	25	y	4	0	
xerces	1.5, 2, 2.5	936	3.0	442	56	31	60	10	y	4	-21	*
log4j	1.0, 1.2, 1.3	1055	1.4	588	30	31	40	29		10	-2	×
xalan	1.0, 1.1	244	1.2	205	32	6	30	6		-2	0	×
	2.4, 2.5, 2.6	2411	2.7	909	38	9	47	9		9	0	×

Fig. 6. Training and test data set properties for Jureczko data, sorted by % defective examples. On the right-hand-side, we show the results from learning. Data is usable if it has a recall of 60% or more and false alarm of 30% or less (and note that, after tuning, there are more usable data sets than before). Results marked with “*” show large improvements in performance, after tuning (lower *pf* or higher *pd*). Data in the three bottom rows, marked with “×”, are performing poorly— that data so few non-defective examples that it is hard for our learners to distinguish between classes.

As mentioned in the last section, this approach depends on having effective predictors for assessing the results. For the Siegmund data, this criteria was relatively easy to achieve. The data in those data sets have a continuous class (runtime of the compiled system) so the performance of a quality predictor can be measured in terms of difference between the predicted run-time *p* of test case items and their actual runtimes *a* using $s = (1 - \frac{abs(a-p)}{a}) \cdot 100\%$ (higher values are better). This paper explores six Siegmund configuration data sets: Berkeley DB (Java and C versions), Apache, SQLite, LLVM, and x264. As a preliminary study, we split that data into equal sized train:test groups and trained a Random Forest Regressor (from the SciKit learn kit [16]) on one half, then applied to the other. After repeating this process 40 times, we achieved nearly perfect accuracy of {98.45%, 99.62%, 94.86%, %, 98.60%, 99.87%% respectively; i.e. we are confident that the predictors from the Siegmund data can assess our plans. (Aside: if the reader doubts that such high scores are achievable, we note that these scores are consistent with those achieved by predictors built by Siegmund et al. [3].)

It proved to be more complicated to commission the Jureczko data sets for this study. For that data, we found that the quality predictors built from this data are far from perfect; However, for some data sets, the predictors could be salvaged using the techniques discussed in this section.

Figure 6 shows our preliminary studies with the Jureczko data. Given access to *V* released versions, we test on version *V* and train on the available data from *V* - 1 earlier releases (as shown in Figure 6, this means that we are training on hundreds to thousands of classes and testing on smaller test suites). Note the three bottom rows marked with ×: these contain predominately defective classes (two-thirds, or more). In such data sets, it is hard to distinguish good from bad (since there are so many bad examples).

In order to identify the presence (or absence) of defects, we can consider using Boolean classes in the Jureczko data (True if defects > 0; False if defects = 0). For such data, quality of the predictor can be measured using

- (1) the probability of detection (a.k.a. “pd” or recall): the percent of faulty classes in the test data detected by the predictor; and
- (2) the probability of false alarm (a.k.a. “pf”): the percent of non-fault classes that are predicted to be defective.

As a preliminary study, we split the Jureczko data into train and test groups. Random Forest (again, from the SciKit learn kit [16]) was built from the training data, then applied to the test data. The “untuned” columns of Figure 6 shows those results. We called a data set “usable” if Random Forest was able to classify majority of the instances correctly. For this purpose, we set a threshold of $pd \geq 60\%$ and $pf \leq 30\%$ to select suitable data sets. With this threshold, however, none of our data sets were suitable for this study.

Fortunately, the “tuned” columns of Figure 6 show that we can salvage some of the data sets. Pelayo and Dick [17] report that defect prediction is improved by SMOTE [18]; i.e. an over-sampling of minority-class examples and an under-sampling of majority-class examples. Also, Fu et al. [19] report that parameter tuning with differential evolution [20] can quickly explore the tuning options of Random Forest to find better settings

for the (e.g.) size of the forest, the termination criteria for tree generation, etc. The rows marked with a \star in Figure 6 show data sets whose performance was improved remarkably by these techniques. For example, in *poi*, the recall increased by 4% while the false alarm rate dropped by 21%. However, as might have been expected, we could not salvage the data sets in the three bottom rows.

In summary, while we cannot trust predictors from some of our Jureczko data sets, we can plan ways to reduce defects in *jedit*, *ivy*, *ant*, *lucene* and *poi*. Accordingly, when this study explores the Jureczko data, we will use these five data sets.

(Aside: One important detail to be stressed here is that, when we applied SMOTE-ing and parameter tunings, those techniques were applied to the training data and *not* the test data; i.e. we took care that no clues from the test set were ever used in this tuning process.)

III. FOUR PLANNING METHODS

This section describes XTREE (which we call Method4) and three alternate methods for learning plans. XTREE, a novel planner introduced in this work, uses the decision tree learner of Figure 7.D. The other methods use the top-down bi-clustering method described in Figure 7.C which recursively divides the data in two using a dimension that captures the greatest variability in the data. We proposed Method 1 and 2 in 2012 [4] while Method 3 comes from research conducted earlier this year [5]. All methods have the properties proposed in §II: (1) they are *local learners* that find plans particular to each test case; and (2) they are *trust-increasing*; i.e. they change examples such that they move *closer* to the training data.

A. Methods

Our description of the methods adopts the following convention. All variables set via our engineering judgement with Greek letters; e.g. α , β , γ , ω . In this paper, we show our current settings to these variables produces useful results. Elsewhere [19], [21], we are exploring tuning methods to find better settings but we have nothing definitive yet to report on auto-tuning planners.

1) Method1 = CD = Centroid Deltas: Summary1: Method1

computes a plan from the difference between where you are (which we will call C_i) and where you want to be (which we will call C_j).

Assumption1: Large data sets can be adequately represented by a few dozen (or so) centroids.

Details1: Method1 clusters project data by reflecting on the independent variables, then reports the delta between the cluster centroids. After clustering training data using the WHERE algorithm of Figure 7.C, Method1 replaces all clusters with a centroid C_i computed from the median value of each continuous/discrete feature. After that, it finds the closest centroid C_j that has a better performance score. For Jureczko data, “better” means fewer defective examples while for the Seigmund data, “better” means lower median runtimes for the examples in that cluster. Method1 then caches the delta between the independent features between C_i and C_j . For continuous features, this delta is $C_j - C_i$. For discrete values, this delta is the value of that feature in C_j . Finally, for every test case, Method1 use the distance measure d shown in Figure 7.B to find the nearest centroid C_i . It then proposes a plan for improving that test case that is the conjunction of all the deltas between C_i and C_j .

2) Method2 = CD + FS = Method1 + Feature Selection : Summary2: Method2 works like Method1 but now the plans only mention the $\beta = 33\%$ most informative features. Hence, Method2’s plans are simpler.

Assumption2: When reasoning about centroids, we can just use features that best distinguish centroids; i.e. whose values appear in just a few centroids. *Details2:* A common result is that the signal in a table of data is mostly contained in a handful of features [22], [23]. Papakroni [24] has tested for this effect in the Jureczko data sets. Papakroni found no loss of efficacy in defect prediction after sorting all features by their information content, then making predictions using (a) all features or (b) just using 33% most informative features.

Based on the above, it might be possible to simplify the plans found by Method1 by pruning back the features in those plans. Following on from Papakroni, our Method2 returns plans containing just the top $\beta = 33\%$ most

informative features. Here, “informative” means that the values of a feature are good for selecting a small set of clusters (ideally, just one). This can be estimated using the Fayyad-Iranni INFOGAIN algorithm [25] of Figure 7.E.

3) *Method3= Best In Cluster (BIC)= Method1 + Gradient: Summary3:* Method3 is like Method1, but it uses more knowledge about the training data.

Assumption3: There exists “gradients” between clusters which, if used, will better guide us to finding beneficial plans. *Details3:* Method3 generates clusters from the training data just like Method1. Following this, it summarizes the clusters C_X using “Best-In-Cluster” B_X , which is (1) the centroid of the cluster for the Jureczko data; (2) the cluster’s member with the fastest runtime for the Seigmund data. Following this, Method3 connects each cluster C_i to a nearest neighbor C_j by a *gradient*. Each gradient has a (bottom,top) end labelled (C_i, C_j) containing the (worst,best) performance scores, respectively.

For each test instance, Method3 finds the nearest gradient, the runs up to the top (best) end C_j , then extracts B_j (which is the best-in-cluster associated with C_j). The returned plan is then computed as a fraction ($\omega = 0.5$) of the distance between the test case and it’s corresponding B_j .

4) *Method4=XTREE= Deltas in Decision Branches: Summary4:* Method4 builds a decision tree, then generates plans from the difference between two branches: (1) the branch to where you are and (2) the branch to where you want to be.

Assumption4: One potential problem with Methods 1,2 and 3 is the *unsupervised* nature of the clustering algorithm (WHERE) that executes without knowledge of the target class. *Supervised* methods assume that it is useful to also reflect on the target class.

Details4: XTREE uses a supervised decision tree algorithm of Figure 7.D to divide the data. Next, XTREE builds plans from the branches of the decision trees using the code of Figure 8. That code asks three questions, the last of which returns the plan:

- 1) What *current* branch does a test case fall in?
- 2) What *desired* branch would the test case want to move to?
- 3) What are the *deltas* between current and desired?

IV. EXPERIMENTAL ANALYSIS

This section describes an experimental design (and results) for evaluating the above four methods.

Figure 7.A: Measuring Variability

For both continuous and discrete values, the variability can be measured using standard deviation (σ) or entropy (e). For continuous values, the standard deviation is given by the formula:

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

where \bar{x} is the mean of the numeric features $x_1, x_2, x_3, \dots, x_n$ and n is the number of data points. For discrete values, the entropy is calculated using the following formula:

$$e = - \sum_{i=1}^n p_i \log_2(p_i)$$

where p_i represents the probability discrete value.

Figure 7.B: Measuring distance

We use Aha et al.’s standard Euclidean distance measure [26]. For F independent features, the measure returns feature (usually set to 1).

$$d(X, Y) = \sqrt{\sum_{i=1}^F w_i \cdot \Delta(X_i, Y_i)^2}$$

Within Δ , if X_i, Y_i are both missing values, then return 1. Otherwise, replace any missing items with values that maximizes the following. For numerics, Δ normalizes X_i, Y_i (to the range 0,1 for min,max) then returns X_i, Y_i . For discrete variables, Δ returns 0,1 if X_i, Y_i are the same,different (respectively).

Figure 7.C: Top-down Clustering with WHERE

WHERE divides data into groups of size $\alpha = \sqrt{N}$. Using this measure, WHERE runs as follows:

- 1) Find two distance cases, X, Y by picking any case W at random, then setting X to its most distant case, then setting Y to the case most distant from X (which requires only $O(2N)$ comparisons).
- 2) Project each case Z to position x on a lines running from X to Y : if a, b are distances Z to X, Y then $x = (a^2 + c^2 - b^2)/(2ac)$.
- 3) Split the data at the median x value of all cases
- 4) For splits larger than $\alpha = \sqrt{N}$, recurse from step1

In terms of related work, the above is similar in approach to Boley's PDDP algorithm [27], but PDDP requires an $O(N^2)$ calculation at each recursive level to find the PCA principle component. Our method, on the other hand, performs the same task with only $O(2N)$ distance calculations using the FASTMAP heuristic [28] shown in step1. Platt [29] notes that FASTMAP is a Nystro'm approximation to the first component of PCA.

Figure 7.D: Top-down division with Decision Trees

Find a split in the values of independent features that most reduces the variability of the dependent feature (measured using Figure 7.A). Construct a standard decision tree using these splits

Figure 7.E: Finding the most informative rows

Discretize all numeric features using the Fayyad-Iranni discretizer [25] (divide numeric columns into bins B_i , each of which select for the fewest cluster ids). Let feature F have bins B_i , each of which contains n_i rows and let each bin B_i have entropy e_i computed from the frequency of clusters seen in that bin (computed from Figure 7.A). Cull the the features as per [24]; i.e. just use the $\beta = 33\%$ most informative features where the value of feature F is $\sum_i e_i \frac{n_i}{N}$ (N is the number of rows)

Using the training data, divide the data using the decision tree algorithm of Figure 7D into groups of size $\alpha = \sqrt{N}$.

For each test item, find the current leaf: take each test instance, run it down to a leaf in the decision tree. After that, find the desired leaf:

- Starting at current, ascend the tree lvl $\in \{0, 1, 2, \dots\}$ levels;
- Identify sibling leaves; i.e. leaf clusters that can be reached from level lvl that are not same as current.
- Using the *score* defined above, find the *better* siblings; i.e. those with a score less than $\gamma = 0.5$ times the mean score of current. If none found, then repeat for lvl+ = 1. Also, return no plan if the new lvl is above the root.
- Return the *closest* better sibling where distance is measured between the mean centroids of that sibling and *current*
- Also, find the *delta*; i.e. the set difference between conditions in the decision tree branch to *desired* and *current*. To find that delta: (1) for discrete attributes, delta is the value from *desired*; (2) for numerics, delta is the numeric difference; (3) for numerics discretized into ranges, delta is a random number selected from the low and high boundaries of the that range.

Finally, return the delta as the plan for improving the test instance.

A. Experimental Design

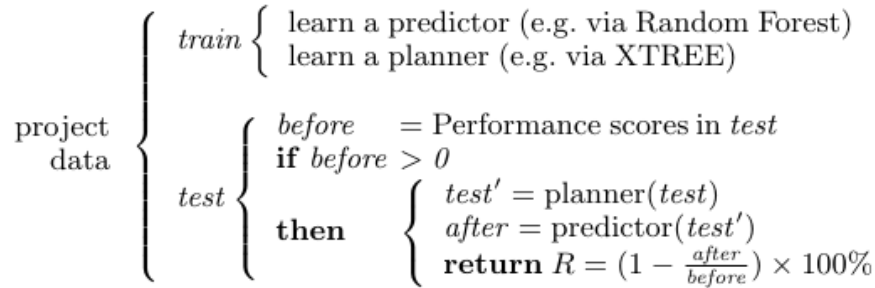


Fig. 9. Experimental design .

1) *A Strategy for Evaluating Planners:* Our experimental design is shown in Figure 9. We divide the project data into two disjoint sets *train* and *test* (so $\text{train} \cap \text{test} = \emptyset$). Next, from the train set, we build both a *planner* and a *predictor*.

Our general framework does not commit to any particular choice of *planner* or *predictor* but, for the purposes of this paper:

- Our *planner* will be one of Methods 1,2,3,4;
- Our *predictor* will be the Random Forest Classifier [30] (for discrete classes) and Random Forest Regressor (for continuous classes) taken from SciKit Learn [16]. We use these data miners since extensive studies have shown these to be amongst the better alternatives for mining software data [31].

As for the *test* data, this is passed to the *predictor* to measure performance statistics related to effectiveness.

If our *predictors* fail to perform effectively on the *test* data, then we cannot trust them to comment on our plans. Accordingly, if that performance is unsatisfactory, we abort. Recall from §?? that this step indicated we should not use some of the Jureczko data.

Else, we (1) apply the *planner* to alter the *test* data; then

(2) apply the *predictor* to the altered data *test'*; then (3) return data on the *before*, *after* predictions expressed as percent improvement, denoted by

$$R = (1 - \frac{\text{after}}{\text{before}}) \times 100\%$$

the following following properties:

- If $R = 0\%$, this means “no change from baseline”;
- If $R > 0\%$, this indicates “improvement over the baseline”;
- If $R < 0\%$, this indicates “optimization failure”.

2) *Statistical Methods:* Our methods use some stochastic algorithms; e.g. WHERE’s selection of “what example to explore first” (see Figure 7.C) and XTREE’ occasional use of a random guess when deciding what part of a discretized range to include in the plan (see Figure 8). Hence, we report the R values seen in 40 repeated runs, with different random number seeds (we use 40 since that is more than the 30 samples needed to satisfy the central limit theorem).

To rank our methods using the results from these 40 repeats, we use the Scott-Knott test recommended by Mittas and Angelis [32].

In accordance to that test, using the median values of each method, we sort a list of $l = 40$ values of R values found in $l_s = 4$ different methods. Then, we split l into sub-lists m, n in order to maximize the expected value of differences in the observed performances before and after divisions. E.g. for lists l, m, n of size l_s, m_s, n_s where $l = m \cup n$:

$$E(\Delta) = \frac{ms}{ls} abs(m.\mu - l.\mu)^2 + \frac{ns}{ls} abs(n.\mu - l.\mu)^2$$

We then apply a statistical hypothesis test H to check if m, n are significantly different (in our case, the conjunction of $A12$ and bootstrapping). If so, Scott-Knott recurses on the splits. In other words, we divide the data if *both* bootstrap sampling and effect size test agree that a division is statistically significant (with a confidence of 99%) and not a small effect ($A12 \geq 0.6$).

For a justification of the use of non-parametric bootstrapping, see Efron & Tibshirani [33, p220-223]. For a justification of the use of effect size tests see Shepperd&MacDonell [34]; Kampenes [35]; and Kocaguenli et al. [36]. These researchers warn that even if a hypothesis test declares two populations to be “significantly” different, then that result is misleading if the “effect size” is very small. Hence, to assess the performance differences we first must rule out small effects using $A12$, a test recently endorsed by Arcuri and Briand at ICSE’11 [37].

3) *Report Format*: Our results are presented in the form of line diagrams like those shown on the right-hand-side of the following example table. The black dot shows the median R value and the horizontal lines stretch from the 25th percentile to the 75th percentile (a region called the inter-quartile range, or IQR).

Example			
Rank	Treatment	Median	IQR
1	XTREE	62	6
2	BIC	50	12
3	CD	44	18
3	CD+FS	43	13

In this example table, the rows are sorted on the median values of each method. Note that all the methods have $R > 0\%$; i.e. all these methods reduced the expected value of the performance score in that experiment while XTREE achieved the greatest reduction (of 62% from the original value).

The above example table has a left-hand-side **Rank** column, computed using the Scott-Knott test described above. In this example table, XTREE is ranked the best, while CD and CD+FS together are ranked the worst. Other Details: and Figure 10 show the effectiveness of our methods seen in 40 repeats with each data set. In these experiments, the dependent variables of Jureczko and Siegmund data sets are discrete and continuous in nature, respectively. Hence, while choosing the predictor, we used Random Forest (1) as a classifier for Jureczko data and (2) as a regressor for Siegmund data. For Siegmund data, we performed a k-fold cross validation. Given the relatively small sample sizes in Apache and BDBJ, we chose $k = 2$. This, we reasoned, would maintain a sufficiently large test data size in order for R to be measured reliably. However, in 2-fold cross validation, the division of data tends to affect the outcome significantly. To counter this, we randomized the order of the data, training on one half while identifying treatment plans on the other, repeating the process 40 times as mentioned above.

The Jureczko data, being temporal in nature, allows us to implement a validation procedure that ensures only past data is ever used to predict future values. Hence, in that data, we used the train/test sets shown in Figure 6. (Aside: note that all the SMOTE-ing and Random Forest tunings (discussed in §??) occurred in the *train* phase of Figure 9).

B. Experimental Results

Recall from our introduction that we are assessing planners on three criteria: *effectiveness*, which is how much they reduce the expected value of the changed examples; *succinctness*, which is how many things we need to change to achieve a plan; and *surprise*, which is how different are the plans from standard truisms.

1) *Effectiveness Results*: Measured in terms of effective-

ness, some data sets were harder to optimize than others. SQL (in Figure 10) defied all our methods for reducing runtimes. XTREE was the only method that could optimize BDBJ (in Figure 10). In general, in most data sets, large reductions were observed:

- An improvement of 60%, as compared to the original baseline in Ant and Lucene of
- An improvement of 77% as compared to the original baseline in BDBC of Figure 10;

Overall, XTREE was most the effective. It was always the top-ranked method with the exception of SQL. Where it ranked better, it had significant improvements in the median performance values (1) In the Jureczko data sets, there was a median improvement of around 10% larger than the next ranked method; (2) In the

Seigmund data sets, improvements as large as 36% greater than the next ranked method were observed.

Rank	Treatment	Median	IQR	Rank	Treatment	Median	IQR
1	XTREE	43	13	1	XTREE	77	8
2	BIC	7	5	2	BIC	66	5
3	CD+FS	0	1	3	CD	-1	1
3	CD	0	1	3	CD+FS	-1	1

Rank	Treatment	Median	IQR	Rank	Treatment	Median	IQR
1	XTREE	28	11	1	XTREE	28	8
2	BIC	5	4	2	BIC	6	2
3	CD	3	5	3	CD	0	0
4	CD+FS	0	1	3	CD+FS	0	1

Rank	Treatment	Median	IQR	Rank	Treatment	Median	IQR
1	XTREE	1	3	1	XTREE	12	1
1	BIC	0	0	2	BIC	2	1
1	CD	0	0	3	CD	0	0
1	CD+FS	-1	0	3	CD+FS	0	0

Fig. 10. Results: Seigmund data sets. Results from 40 repeats. Ratios of (1) software runtimes expected in the test examples after alteration by a planner to (2) the sum of the software runtimes in the original test set. Larger median values are better.

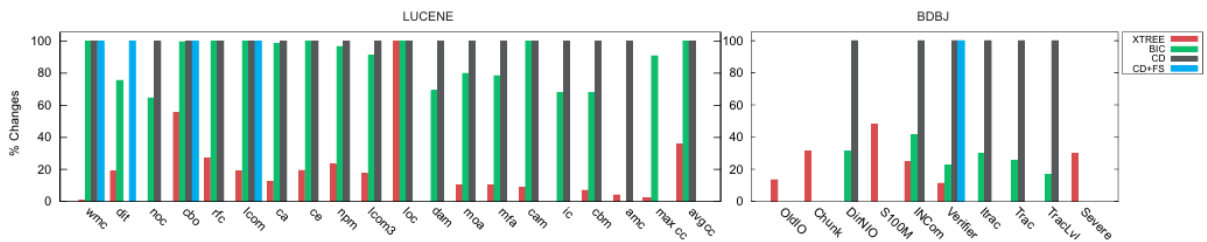


Fig. 11. Percent frequency for how often certain feature was changed by a plan.

2) *Succinctness Results:* Figure 11 reports the percent of times in the 40 repeats that a method proposed changing a feature. The left-hand-side plot of that figure reports results from one of the Jureczko data sets (*lucene*) and the right-hand-side shows a Siegmund data set (*BDBJ*).

In these plots, *more* succinct a planning method, *fewer* the frequency (in percent) where it recommends changing a particular feature (i.e. the vertical bars in that plot are *lower*). For example, XTREE’s plans were usually succinct – in all data sets, XTREE changes around a fifth of the features (see Figure 12). On other hand, Method1 (CD) was the least succinct since it wanted to change all features (observe the change frequencies as high as 100% for all features). Method1’s policy of “change everything” might be acceptable if this approach lead to the most effective changes. However, in and Figure 10, there is no evidence for this.

An interesting feature of Figure 11 was that fewer things were changed in the Siegmund data sets like *BDBJ* than in the Jureczko data set like *lucene*. It turns out that this holds true across nearly all our data sets. Figure 12 summarizes all the change frequencies for all data sets. As with Figure 11, there are fewer changed features in the Siegmund data than in the defect prediction data. One explanation for that is the nature of the features: the Jureczko data sets have continuous features while the Siegmund data has binary independent features (where settings can be turned “ON” or “OFF”). While the features of the Jureczko data sets are not constrained, those of Siegmund data sets are subject to several constrains (dictated by their feature models, see Figure 5). It is possible that planners in Method1, Method2, and Method3 find very few valid settings that can improve the performance scores, thereby making insignificant changes. As a result all methods, except XTREE, fail to optimize the runtime, as evidenced by the results in Figure 10.

		XTREE	BIC	CD	CD+FS
Jureczko data	Ant	20%	87%	80%	10%
	Ivy	21%	81%	70%	20%
	Jedit	19%	92%	90%	20%
	Lucene	19%	86%	95%	25%
	Poi	13%	79%	85%	20%
	mean:	18%	85%	84%	19%
Siegmund data	Apache	27%	5%	33%	22%
	BDBJ	7%	8%	3%	15%
	LLVM	33%	4%	36%	9%
	X264	24%	12%	18%	12%
	BDBC	25%	9%	16%	5%
	SQL	10%	6%	10%	23%
	mean:	21%	7%	19%	14%

Fig. 12. Average number of features whose values are changed by a planner.

3) *Surprising Results*: If a planner only ever reported conclusions that were already known, then that planner offers little value over “just use established wisdom”. Accordingly, we studied our results for plans that were somewhat counter- intuitive.

Such a surprising plan can be seen in *lucene*. Recall the standard advice for OO systems: build classes that are internally cohesive with low coupling to other parts of the system [15]. We can assess the relevance of this advice to specific projects by checking how often a planner changes the coupling-related features:

- *ca*: afferent couplings = # classes using this class;
- *ce*: efferent couplings = # classes used by this class.
- *cbm*: coupling between methods = # new/redefined methods to which all the inherited methods are coupled
- *cbo*: coupling between objects = a value that increases when the methods of one class access services of another.
- *ic*: inheritance coupling = # parent classes which a given class is coupled (including methods and variables inherited)

In many of our results with the Jureczko data, it was indeed true that the changes proposed by XTREE lead to lower coupling. However, the *lucene* results were quite different and rather surprising.

In the *lucene* XTREE results of Figure 11, the most frequent change was to alter the lines of code in a class (see the tallest red histogram in that figure on the *loc*, or lines of code). Looking at the logs of our planner, we can see that XTREE’s proposed change is to *reduce* the size of a class. The only way to do that, while keeping the same functionality, is to create a network of smaller classes that interact to produce that functionality. That is, we would need to *increase* the coupling of those classes to achieve XTREE’s plan.

In theory, increasing coupling between classes complicates and confuses a class design. But the *lucene* XTREE results of Figure 11 rarely proposes changing the coupling features *ca*, *ce*, *cbm*, *ic* (in fact, XTREE never proposes any change to *ic*). The only coupling issue that XTREE usually adds to its plans is *cbo* (which appears 55% of the time in Figure 11). But note that this is *object* coupling measure, not class coupling. So here XTREE is warning against, say, some factory class generating a large community of agents, all of the same class, who coordinate on some task. This is a different issue to the class redesign issue that would be triggered by altering *loc*.

In summary, XTREE satisfies that criteria that, sometimes, it produces surprising plans. At least for the *lucene* data set, we can see advice that recommends *increasing coupling* to reduce defects.

V. THREATS TO VALIDITY

As with any empirical study, biases can affect the final results. Therefore, any conclusions made from this work must be considered with the following issues in mind.

A. Learning Bias

For building the defect predictors in this study, we elected to use Random Forest and Random Forest Regressors. We chose this approach, based on its reputation for having the better performance of 21 other learners for defect prediction [31]. Data mining is a large and active field and any single study can only use a small subset of the known classification algorithms. That said, we have taken care to document in this paper the decisions made by engineering judgement that could effect our conclusions. The above code used a set of variables which future work should vary in order to test the internal validity of our conclusions:

- All our planners divide data into groups of size $a = \sqrt{N}$
- Method2 used the top $\beta = 33\%$ most informative features (ranked using INFOGAIN);
- Method3 uses $\omega = 0.5$ to find its deltas;
- XTREE considered a sibling useful if it’s score was less than $\gamma = 0.5$ times the mean score of the current leaf. ✓

B. Sampling bias

Sampling bias threatens any data mining experiment; i.e., what matters there may not be true here. For example, the data sets used here comes from two sources (Seigmund et al. and Jureczko et al.) and any biases in their selection procedures threaten the validity of these results. That said, the best we can do is define our methods and publicize our data and code so that other researchers can try to repeat our results and, perhaps, point out a previously unknown bias in our analysis. Hopefully, other researchers will emulate our methods in order to repeat,

refute, or improve our results.

C. Evaluation Bias

Another threat to validity of this work is our use of predictors learned on the training data to assess the impact of our planners. This issue was discussed in detail in §II-C. To those notes, we add a few more details. If possible, planners should be assessed via some external oracle that can accurately assess new examples. For example, in search-based software engineering, examples can be assigned objective scores via some model. In this approach, a changed example can be assessed by generating actual objective scores from the model

Rank	Treatment	Median	IQR	
1	XTREE	59	9	•
2	BIC	5	1	•
2	CD+FS	-7	100	—•—
2	CD	-9	77	—•—

Fig. 13. Methods 1,2,3,4 applied to some ground-truth data (in this case, the POM3 model). Values collected from 40 repeated runs of each method with different random seeds. Results show the efficacy of XTREE in reducing the total overall cost in the original test data, when other planners fail to do so.

The POM3 model [38], [39] is a tool for exploring management challenges. POM3 implements the Boehm and Turner model of agile programming [40] where teams select tasks as they appear in the scrum backlog. POM3 can study the implications of different ways to adjust task lists in the face of shifting priorities. The model outputs estimated task completion rates; programmer idle rates; and total overall cost. POM3 models requirements as a dependency tree. A single requirement in the tree of a prioritization value has a cost, along with a list of child-requirements and dependencies. Before any requirement can be satisfied, its children and dependencies must first be satisfied. POM3 simulates changing priorities by making teams aware of random items in the requirements tree at random intervals, thus forcing teams to constantly readjust their “to do” lists. For further details on this model see [38], [39], and [40].

Figure 13 shows results from using POM3 as an oracle to assess our planning methods and their ability to reduce the project cost. In this experiment, we generated a training and testing set with 1000 randomly generated instances, which we passed to our four methods. 40 times, we let those methods propose changes to those projects. For assessment purposes, the changed projects were then fed back to the POM3 oracle.

Using this approach, it is possible to assess the value of a plan by measuring its effectiveness with respect to some ground truth (in this case, the POM3 model). As shown in Figure 13, XTREE reduces the cost by 59% thereby passing this assessment. These results from the POM3 model further endorse our previous conclusions; i.e. compared to three other methods, XTREE’s supervised methods are best for generating plans on how to change example projects.

VI. RELATED WORK

A. Planning in AI

The XTREE planner is somewhat different to the logic-based planners explored by classical AI. Those kinds of planners employ a logical procedure [41] that seeks an ordering on *operators* to take some domain *state* from a *start* state to a *goal* state. This classical logical approach is known to suffer from computational bottlenecks [42]. On the other hand, tools like XTREE will scale to any domain that can generate decision trees.

B. Evaluating Changes

Some organizations have the resources to run repeated trials to assess project changes. For example, in one study, Bente et al. reported results where the same specification was developed by four different organizations [43]. Given those kind of resources, it would be possible to (say) take a code base, assign it to different teams, make these teams adopt different polices, then check in 12 months time which teams have fewer defects than the others.

Seldom do industrial or research groups have access to the kinds of resources needed for this kind study (evidence: in the six years since the publication of that work, we know of only one similar study to Bente et al.). Also, given the diversity of modern software projects, it might be unreasonable to demand that all proposed changes for all projects are always evaluated by something like the Bente et al. study. Hence, this paper has used data miners to build an oracle that can assess changed examples. The advantage of our approach is that it required far less resources to assess the effectiveness of proposed changes to a project.

C. Search based SE

Another way to propose changes to software artifacts is via some search-based method [44], [45]. Such SBSE methods are evolutionary programs that make extensive changes to some initial sample of project data (perhaps 100s to 100,000s of mutations). Each of these mutations is reassessed using some domain model. Examples of these algorithms include GALE, NSGA-II, NSGA-III, SPEA2, IBEA, particle swarm optimization, MOEA/D, etc. [21], [46]–[51].

One problem with these SBSE methods is that they can make extensive mutations to the data they are exploring. In the language of §II-C, these methods may not be *trust-increasing* since those algorithms make no attempt to prevent new examples from mutating away from the kinds of data used to commission the model (in which case, we would start doubting the model's output).

Another issue with standard search-based SE methods is that they require ready access to trustworthy domain model that can offer an assessment of newly generated examples. While some domains have such models (e.g. see the COCOMO effort estimation model used in the last section), our experience is that many others do not. For example, consider software defect prediction and all the intricate issues that may lead to defects in a product. A model that includes *all* those potential issues would be very large and complex. Further, the empirical data required to validate any/all parts of that model can be hard to find.

What we would recommend is a two-pronged policy. In domains with ready access to trusted models, we recommend the kinds of tools that are widely used in the search-based software engineering community such as GALE, NSGA-II, NSGA-III, SPEA2, IBEA, particle swarm optimization, MOEA/D, etc. [21], [46]–[51]. Otherwise, we recommend tools like XTREE.

VII. CONCLUSION

The planner proposed in this paper proposes changes to software project details in order to improve the expected value of the performance scores of that part of the project. To evaluate these planners, data miners can be used to build oracles to assess planners. Such planners should be *trust-increasing*; i.e. they should propose changes that generate changed examples that are closer to the training data of the data miner. One caveat here is that the evaluations we can make on the planner are only as good as the predictive performance of the data miner. Hence, if domain data does not support satisfactory predictors, then planning in that domain cannot be evaluated. Four planners were assessed here for the tasks of reducing defects and runtimes. Three of the methods come from our prior publications [4], [5], and the conclusion of this paper is that a novel fourth method clearly outperforms the other three (measured in terms of effectiveness, succinctness, and surprise). We conjecture that XTREE worked better than the rest since:

- It uses a supervised method to divide the data and;
- When planning how to move examples to better classes, it is best to reflect over differences between those classes.

VIII. FUTURE WORK

Future work in this area could explore numerous questions. For example, XTREE, as used here, sought improvements in a one goal (the class variable). Does XTREE work for multi-goal reasoning?

Also, the XTREE algorithm seems quite general to any data of table with rows containing weighted classes (so we can distinguish “bad” rows from “better” ones). Does XTREE works on domains (other than the defect/runtime data explored here)?

As an example of a domain that might benefit from XTREE, recent results raise doubts about the value of changing code to remove “bad smells” [52]. Can XTREE be used as a “bad smell” detector to select the subset of possible refactorings that have the most potential benefit?

As to scalability, XTREE is a post-processor to a decision tree algorithm. Hence, in theory, XTREE works on domain where data miners can generate decision trees. Given the current state of the art in Big Data, can XTREE be applied to very large data sets?

We discussed above in §II-B the general conclusions of Passos, Jørgensen et.al [6], [7]; i.e. software developers are reluctant to surrender their old biases in the face of new data. Accordingly, it must be asked if the *mental resistance* of developers will prevent them applying XTREE's automatically generated recommendations of tools? Note this this issue is not just a concern for XTREE, but also for any automatic tool proposing refactorings.

There are many more methods for generating plans and no one paper can survey them all. For example, this paper has not explored variations to the α , β , and γ parameters that controlled XTREE. Would we get better results if we varied those parameters?

That said, the goal of this paper was not to claim that (e.g.) XTREE is some absolute optimal algorithm. Rather, it is was to offer a baseline result (with XTREE) and an evaluation strategy that can assess if alternate methods are better than XTREE. The authors of this paper would actively support other teams exploring this method (with or without using our current code base).

REFERENCES

- [1] T. Menzies and T. Zimmermann, "Goldfish bowl panel: Software development analytics," in ICSE'12, pp. 1032–1033.
- [2] M. Jureczko and L. Madeyski, "Towards identifying software project clusters with regard to defect prediction," in Proceedings of the 6th International Conference on Predictive Models in Software Engineering, ser. PROMISE '10. ACM, 2010, pp. 9:1-9:10
- [3] N. Siegmund, S. S. Kolesnikov, C. Kästner, S. Apel, D. Batory, M. Rosenmüller, and G. Saake, "Predicting performance via automated feature-interaction detection," in Proceedings of the 34th International Conference on Software Engineering. IEEE Press, 2012, pp. 167–177.
- [4] R. Borges and T. Menzies, "Learning to Change Projects," in Proceedings of PROMISE'12, Lund, Sweden, 2012.
- [5] R. Krishna and T. Menzies, "Actionable = cluster + constraint?" ACTION'15: an ASE'15 workshop, 2015.
- [6] C. Passos, A. P. Braun, D. S. Cruzes, and M. Mendonca, "Analyzing the impact of beliefs in software project practices," in ESEM'11, 2011.
- [7] M. Jørgensen and T. M. Gruschke, "The impact of lessons-learned sessions on effort estimation and uncertainty assessments," Software Engineering, IEEE Transactions on, vol. 35, no. 3, pp. 368–383, May-June 2009.
- [8] N. Bettenburg, M. Nagappan, and A. E. Hassan, "Towards improving statistical modeling of software engineering data: think locally, act globally!" Empirical Software Engineering, pp. 1–42, 2014. [Online]. Available: <http://dx.doi.org/10.1007/s10664-013-9292-6>
- [9] 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, 2011.
- [10] Y. Yang, Z. He, K. Mao, Q. Li, V. Nguyen, B. W. Boehm, and R. Valerdi, "Analyzing and handling local bias for calibrating parametric cost estimation models," Information & Software Technology, vol. 55, no. 8, pp. 1496–1511, 2013.
- [11] L. L. Minku and X. Yao, "Ensembles and locality: Insight on improving software effort estimation," Information & Software Technology, vol. 55, no. 8, pp. 1512–1528, 2013.
- [12] T. Menzies, A. Butcher, D. R. Cok, A. Marcus, L. Layman, F. Shull, B. Turhan, and T. Zimmermann, "Local versus global lessons for defect prediction and effort estimation," IEEE Trans. Software Eng., vol. 39, no. 6, pp. 822–834, 2013, available from <http://menzies.us/pdf/12localb.pdf>.
- [13] 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>.
- [14] D. Posnett, V. Filkov, and P. Devanbu, "Ecological inference in empirical software engineering," in Proceedings of ASE'11, 2011.
- [15] H. Dhama, "Quantitative models of cohesion and coupling in software," Journal of Systems and Software, vol. 29, no. 1, pp. 65–74, 1995, oregon Metric Workshop. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/016412129400128A>
- [16] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine Learning in Python," Journal Machine Learning, vol. 12, pp. 2825–2830, 2012. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2078195> <http://arxiv.org/abs/1201.0490>

- [17] L. Pelayo and S. Dick, "Applying novel resampling strategies to software defect prediction," in *Fuzzy Information Processing Society, 2007. NAFIPS '07. Annual Meeting of the North American*, June 2007, pp. 69–72.
- [18] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [19] W. Fu, T. Menzies, and X. Shen, "Tuning for software analytics: is it really necessary?" Submitted to *ICSE '16*, 2016.
- [20] R. Storn and K. Price, "Differential Evolution — A Simple and Efficient Heuristic for global Optimization over Continuous Spaces," *Journal of Global Optimization*, vol. 11, no. 4, pp. 341–359, 1997.
- [21] J. Krall and T. Menzies, "GALE: Geometric Active Learning for Search-based Software Engineering," *IEEE Transactions on Software Engineering (to appear)*, 2015.
- [22] M. A. Hall and G. Holmes, "Benchmarking attribute selection techniques for discrete class data mining," *Knowledge and Data Engineering, IEEE Transactions on*, vol. 15, no. 6, pp. 1437–1447, 2003.
- [23] R. Kohavi and G. H. John, "Wrappers for Feature Subset Selection," *Artificial Intelligence*, vol. 97, no. 1-2, pp. 273–324, 1997. [Online]. Available: citeseer.nj.nec.com/kohavi96wrappers.html
- [24] V. Papakroni, "Data carving: Identifying and removing irrelevancies in the data," Master's thesis, Lane Department of Computer Science and Electrical Engineering, West Virginia University, 2013.
- [25] U. M. Fayyad and I. H. Irani, "Multi-interval Discretization of Continuous-valued Attributes for Classification Learning," in *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence*, 1993, pp. 1022–1027.
- [26] D. W. Aha, D. Kibler, and M. K. Albert, "Instance-based learning algorithms," *Machine learning*, vol. 6, no. 1, pp. 37–66, 1991.
- [27] D. Boley, "Principal direction divisive partitioning," *Data Min. Knowl. Discov.*, vol. 2, no. 4, pp. 325–344, December 1998.
- [28] 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*, 1995, pp. 163–174.
- [29] J. C. Platt, "FastMap, MetricMap, and Landmark MDS are all Nystro^m algorithms," in *In Proceedings of 10th International Workshop on Artificial Intelligence and Statistics*, 2005, pp. 261–268.
- [30] L. Breiman, "Random forests," *Machine learning*, pp. 5–32, 2001. [Online]. Available: <http://link.springer.com/article/10.1023/A:1010933404324>
- [31] S. Lessmann, B. Baesens, C. Mues, and S. Pietsch, "Benchmarking classification models for software defect prediction: A proposed framework and novel findings," *Software Engineering, IEEE Transactions on*, vol. 34, no. 4, pp. 485–496, July 2008.
- [32] N. Mittas and L. Angelis, "Ranking and clustering software cost estimation models through a multiple comparisons algorithm," *IEEE Trans. Software Eng.*, vol. 39, no. 4, pp. 537–551, 2013.
- [33] B. Efron and R. J. Tibshirani, *An introduction to the bootstrap*, ser. *Mono. Stat. Appl. Probab.* London: Chapman and Hall, 1993.
- [34] M. J. Shepperd and S. G. MacDonell, "Evaluating prediction systems in software project estimation," *Information & Software Technology*, vol. 54, no. 8, pp. 820–827, 2012.
- [35] V. B. Kampenes, T. Dyba, J. E. Hannay, and D. I. K. Sjøberg, "A systematic review of effect size in software engineering experiments," *Information & Software Technology*, vol. 49, no. 11-12, pp. 1073–1086, 2007.
- [36] E. Kocaguneli, T. Zimmermann, C. Bird, N. Nagappan, and T. Menzies, "Distributed development considered harmful?" in *Proceedings - International Conference on Software Engineering*, 2013, pp. 882–890.
- [37] A. Arcuri and L. Briand, "A practical guide for using statistical tests to assess randomized algorithms in software engineering," in *ICSE'11*, 2011, pp. 1–10.
- [38] B. Boehm and R. Turner, "Using risk to balance agile and plan-driven methods," *Computer*, no. 6, pp. 57–66, 2003.

- [39]D. Port, A. Olkov, and T. Menzies, "Using simulation to investigate requirements prioritization strategies," in Proceedings of the 2008 23rd IEEE/ACM International Conference on Automated Software Engineering. IEEE Computer Society, 2008, pp. 268–277.
- [40]B. Boehm and R. Turner, *Balancing agility and discipline: A guide for the perplexed*. Addison-Wesley Professional, 2003.
- [41]R. E. Fikes and N. J. Nilsson, "Strips: A new approach to the application of theorem proving to problem solving," pp. 189–208, 1971.
- [42]T. Bylander, "The computational complexity of propositional STRIPS planning," pp. 165–204, 1994.
- [43]B. C. D. Anda, D. I. K. Sjøberg, and A. Mockus, "Variability and reproducibility in software engineering: A study of four companies that developed the same system," *IEEE Transactions on Software Engineering*, vol. 35, no. 3, pp. 407–429, 2009.
- [44]M. Harman, S. Mansouri, and Y. Zhang, "Search Based Software Engineering: A Comprehensive Analysis and Review of Trends Techniques and Applications," *Engineering*, no. TR-09-03, pp. 1–78, 2009. [Online]. Available: <http://discovery.ucl.ac.uk/170689/>
- [45]M. Harman, P. McMinn, J. De Souza, and S. Yoo, "Search based software engineering: Techniques, taxonomy, tutorial," *Search*, vol. 2012, pp. 1–59, 2011. [Online]. Available: <http://discovery.ucl.ac.uk/1340709/>
- [46]K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A Fast Elitist Multi-Objective Genetic Algorithm: NSGA-II," *IEEE Transactions on Evolutionary Computation*, vol. 6, pp. 182–197, 2002.
- [47]E. Zitzler, M. Laumanns, and L. Thiele, "SPEA2: Improving the Strength Pareto Evolutionary Algorithm for Multiobjective Optimization," in *Evolutionary Methods for Design, Optimisation, and Control*. CIMNE, Barcelona, Spain, 2002, pp. 95–100.
- [48]E. Zitzler and S. Künzli, "Indicator-Based Selection in Multiobjective Search," in *Parallel Problem Solving from Nature - PPSN VIII*, ser. Lecture Notes in Computer Science, X. Yao, E. Burke, J. Lozano, J. Smith, J. Merelo-Guervos, J. Bullinaria, J. Rowe, P. Tinó, A. Kabań, and H.-P. Schwefel, Eds. Springer Berlin Heidelberg, 2004, vol. 3242, pp. 832–842. [Online]. Available: http://dx.doi.org/10.1007/978-3-540-30217-9_84.
- [49]K. Deb and H. Jain, "An Evolutionary Many-Objective Optimization Algorithm Using Reference-Point-Based Nondominated Sorting Approach, Part I: Solving Problems With Box Constraints," *Evolutionary Computation, IEEE Transactions on*, vol. 18, no. 4, pp. 577–601, Aug. 2014.
- [50]X. Cui, T. Potok, and P. Palathingal, "Document clustering using particle swarm optimization," . . . *Intelligence Symposium*, 2005. . . ., 2005. [Online]. Available: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=1501621
- [51]Q. Zhang and H. Li, "MOEA/D: A Multiobjective Evolutionary Algorithm Based on Decomposition," *Evolutionary Computation, IEEE Transactions on*, vol. 11, no. 6, pp. 712–731, Dec. 2007.
- [52]D. Sjøberg, A. Yamashita, B. Anda, A. Mockus, and T. Dyba, "Quantifying the effect of code smells on maintenance effort," *Software Engineering, IEEE Transactions on*, vol. 39, no. 8, pp. 1144–1156, Aug. 2013.