

An Empirical Analysis on Web Service Anti-Pattern Detection using a Machine Learning Framework

Lov Kumar
NIT Rourkela, India
lovkumar505@gmail.com

Ashish Sureka
Ashoka University, India
ashish.sureka@ashoka.edu.in

Abstract—Web Services are application components characterised by interoperability, extensibility, distributed application development and service oriented architecture. A complex distributed application can be developed by combing several third-party web-services. Anti-patterns are counter-productive and poor design and practices. Web-services suffers from a multitude of anti-patterns such as God object Web service and Fine grained Web service. Our work is motivated by the need to build techniques for automatically detecting common web-services anti-patterns by doing a static analysis of the source code implementing a web-service. Our approach is based on the premise that summary values of object oriented source code metrics computed at a web-service level can be used as a predictor for anti-patterns. We present an empirical analysis of 4 data sampling technique to encounter the class imbalance problem, 5 feature ranking technique to identify the most informative and relevant features and 8 machine learning algorithms for predicting 5 different types of anti-patterns on 226 real-world web-services across several domains. We conclude that it is possible to predict anti-patterns using source code metrics and a machine learning framework. Our analysis reveals that the best performing classification algorithm is Random Forest, best performing data sampling technique is SMOTE and the best performing feature ranking method is OneR.

Index Terms—Anti-Patterns, Empirical Software Engineering, Feature Selection, Imbalanced Learning, Machine Learning, Predictive Modeling, Service Oriented Architecture, Software Analytics, Source Code Analysis, Source Code Metrics, Web Services.

I. RESEARCH MOTIVATION AND AIM

Web Services are used for building distributed software systems based on the Service Oriented Architecture (SOA) paradigm [1][2][3]. A complex distributed system can be developed by combining web-services from third-party providers. A web-service providing a certain functionality defined by the web service interface and description is published by a web-service provider which can be discovered and invoked by a distributed application [1][2][3]. Web services are developed based on pre-defined standards to support interoperability. Similar to object oriented systems, web services also suffers from anti-patterns due to bad programming practises, design and implementation [4][5][6][7][8][9][10]. Anti-patterns in web-services leads to maintenance and evolution related problems. There are several web-services anti-patterns [4][5][6][7][8][9][10]. The five web-services anti-patterns which we consider in the study presented this paper are: GOWS : God object Web service, FGWS : Fine grained

Web service, DWS : Data Web service, CWS : Chatty Web service, AWS: Ambiguous Web service. For example, the GOWS anti-pattern results because of implementing a lot of methods related to multitude of business and technical abstractions into a single service. AWS anti-pattern has implications in reusability and discoverability of web-services due to ambiguous names. Previous research shows that anti-patterns in web-services is a serious and important problem and there are lack of approaches for automatically detecting web-service anti-patterns from their source code and Web Service Definition Language (WSDL) documents [6][7][8][9][10]. For example, Ouni et al. mention that automatic detection of anti-patterns in web services is relatively unexplored [6]. Their analysis reveals that there are lack of approaches for anti-pattern detection in web-services and research in this direction is still in its infancy [6]. The first approach for web-service anti-pattern detection was proposed by Ouni et al. in the year 2015 [6]. Similarly, Palma et al. also mention that extensive usage and importance of web-services, there are no or very less work on automated approaches for web-services anti-pattern detection [8]. They emphasize the need for building automated approaches for web-services anti-pattern identification to improve web-services maintenance and evolution [8].

The work presented in this paper is motivated by the need to build tools and techniques for automatic detection of anti-patterns in web-services. In particular, our motivation is to investigate the application of machine learning based techniques for building predictive models from source code metrics as features for the task of detecting web-services anti-patterns. The specific research aims of the work presented in this paper are the following:

- 1) To investigate the application of different machine learning classification algorithms, feature selection and ranking techniques and data sampling techniques for the task of web-services anti-pattern detection using source code metrics as features. To examine the correlation between several object oriented metrics derived from the source code implementing a web-service and presence of common web-services anti-patterns.
- 2) To conduct an in-depth empirical analysis on real-world web-services dataset and investigate the performance of the machine learning algorithms and object oriented source code metrics based features for predicting web-

services anti-patterns. To examine the relative performance of various classification algorithms, feature ranking techniques, data sampling methods to encounter the class imbalance problem and object oriented source code metrics using Area Under an ROC Curve (AUC) and statistical hypothesis based testing approaches.

II. RELATED WORK AND RESEARCH CONTRIBUTIONS

Ouni et al. propose a search based approach for web-service anti-pattern detection [6][7]. Their approach is based on Genetic Programming (GP) which generates anti-pattern identification rules from examples of web-service anti-patterns [6][7]. Mateos et al. conduct a research on web-service anti-patterns and conclude that there is a correlation between object-oriented metrics in the source code implementing a web-service and presence of anti-patterns in the Web Service Description Language (WSDL) of a web-service [11]. Their research reveals several bad practises in public WSDL documents which gets manifested in the automatically generated WSDL documents because of poor implementation practises in the source code [11]. Their work consists of investigating the feasibility of web-service anti-patterns using object oriented metrics [11]. Coscia et al. present an empirical analysis on the influence of object oriented source code metrics driven code refactoring on the WSDL file and document quality [12]. Their work is based on the premise that there is a correlation between the occurrences of web-service anti-patterns manifested in WSDL and several source code metrics implementing the web service [12].

Palma et al. propose an approach for detecting anti-patterns in web-services [8][9]. They specify 10 web-services anti-pattern and conduct experiments belonging to the weather and financial domain. Their detection algorithm is based on identifying the relevant properties of web-services specific anti-patterns [8][9]. Moha et al. present a method for

web-services anti-pattern detection which is able to identify anti-patterns such as Multi Service and Tiny Service. They conduct experiments on a service based system in the domain of home automation are able to achieve a precision of around 90% [10]. Nayrolles et al. mention that there is no tool support which exists for the detection of anti-patterns in web-services and service oriented systems [13]. They describe a prototype tool called as SODA which is based on specifying anti-patterns in the form of rule cards and then generating detection algorithms which conforms to the web-services anti-pattern specifications [13].

Research Contributions : In context to existing work, the study presented in this paper makes two main research contributions. The first novel and unique research contribution of our work is the application of 18 object oriented source code metrics as features or independent variables, 8 classification algorithms, 5 feature selection and ranking technique and 4 data sampling techniques to counter the class imbalance problem for the task of detecting 5 different types of web-service anti-patterns. The second contribution of our work to the body of knowledge on web-service anti-pattern detection is an in-depth empirical analysis of the proposed techniques and frameworks on a publicly available dataset consisting of 226 web-services across several domains. We conduct a series of experiments investigating the performance of several machine learning parameter combinations and examine their relative effectiveness using AUC based performance metrics and statistical hypothesis based testing.

III. SOLUTION APPROACH AND RESEARCH FRAMEWORK

Figure 1 illustrates our research framework and solution approach. As shown in Figure 1, our approach consists of multiple steps. The first step consists of computing the features of attributes from the 226 web-services. The features

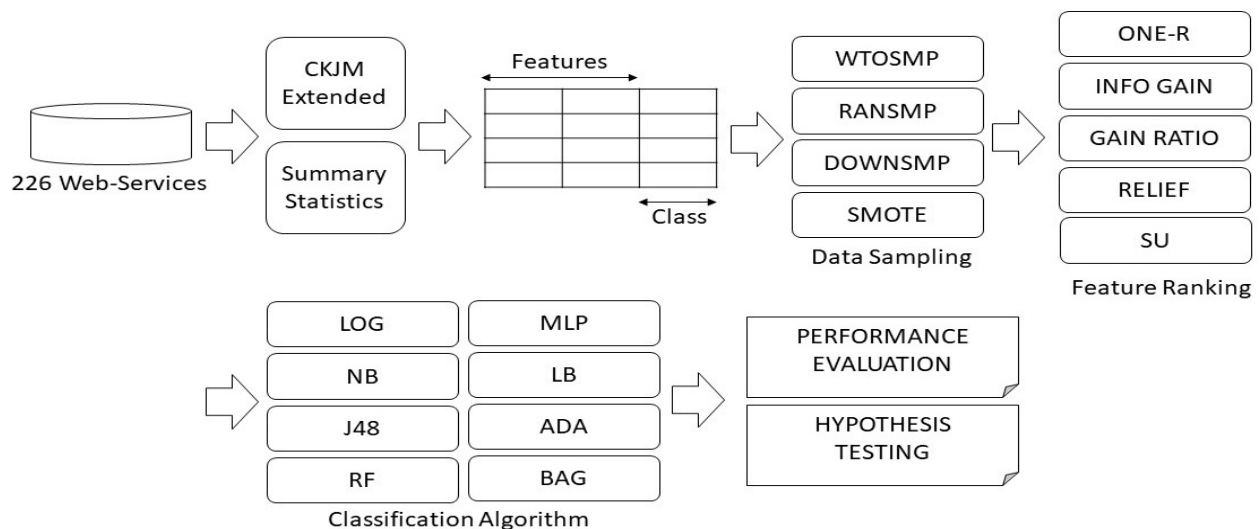


Fig. 1: Research Framework and Solution Approach

are computed by obtaining the object oriented metrics using the CKJM Extended Tool and then computing summary statistics from it. The dataset has imbalanced class distribution and hence we apply several data sampling techniques to encounter the data imbalance problem. We apply several feature ranking techniques to identify informative features and remove irrelevant or redundant features (refer to Figure 1). As shown in Figure 1, we apply different classification algorithms to build predictive models and perform a comparison between the classifiers. Finally, we evaluate the performance of the classifiers using performance evaluation metrics and draw conclusions based on the results of the statistical hypothesis testing.

Source Code Metrics: We compute the object oriented source code metrics using the CKJM extended tool¹. The object oriented metrics are computed at the web-service level for the Java source code implementing the web-services. The CKJM extended tool computes several Chidamber and Kemerer (defined in [14][15]) Java Metrics. The metrics are computed by processing the bytecode of the compiled Java files. We compute a total of 18 size and structure metrics such as WMC: Weighted methods per class, DIT: Depth of Inheritance Tree, NOC: Number of Children, CBO: Coupling between object classes, LCOM3: Lack of cohesion in methods Henderson-Sellers version, LCO: Lines of Code, DAM: Data Access Metric. A complete list of the metrics with their abbreviations and expanded-form is provided in the CKJM tool documentation. We compute the aggregate and summary values at the web-service level once the 18 object oriented metrics values are computed for every class in the web service. The values which we compute at the web service level are: minimum, maximum, mean, median, first quartile, third quartile, standard deviation, skewness and kurtosis. We compute these values for all the 18 metrics except in cases where such values cannot be computed. These summary values obtained from the data are the attributes for the classification problem.

Data Sampling Techniques to Learn from Imbalanced Data: Our experimental dataset has an imbalanced class distribution. The number of web services containing a particular anti-pattern is significantly lower (5% to 10% of the dataset) than the number of web services without the anti-pattern. Dealing with an imbalanced data containing unequal instanced of different classes is one the technical challenges of our problem. We use three different types of sampling techniques to overcome the unbalanced 2-class machine learning problem. We apply a technique called as Downsampling (DOWNSMP) which creates a balanced dataset such that the number of instanced of the minority class is the same as the number of instances of the majority class [16][17]. Another sapling technique which we use is RANSMP which is a nave oversampling approach. RANSMP

consists of generating new and more samples for the class which is under-represented and is in minority. We apply a technique called as SMOTE (Synthetic Minority Over-sampling Technique) which is based on the combination of over-sampling the minority class and over-sampling the majority class [16].

Feature Ranking and Selection: Liu et al. present a survey of several feature selection algorithms [18]. We apply five different feature selection approaches (refer to Figure 1): OneR - FR1, Information Gain (IG) - FR2, Gain Ratio (GR) - FR3, RELEIF (RF) - FR4 and Symmetric Uncertainty (SU) - FR5. We also consider a case with no feature ranking (AM : All Metrics).

Learning Algorithms: As illustrated in Figure 1, we apply eight different types of classifiers or learning algorithms: J48 : Decision Tree, RF : Random Forest, LB : Logit Boost, ADA: AdaBoost, BAG : Bagging, MLP : Multilayer Perceptron, NB : Naive Bayes, LOG: Logistic regression. Lim et al. present a comparison of the predictive accuracy, complexity and training time of 33 learning algorithms [19]. We apply both classical algorithms (such as Decision Tree and Nave Bayes) as well as relatively modern and new statistical algorithms to examine the generalizability of our empirical results and findings and identify the best performing classification algorithm.

Performance Evaluation : Effectiveness of a learning algorithm or classifier can be measured using several different performance metrics. Sokolova present a detailed comprehensive analysis of the performance measures for predictive modelling and classification task [20]. We use metrics such as Area under the ROC (Receiver Operating Characteristics) Curve (AUC), f-measure and accuracy to measure the performance of the learning algorithm. We use f-measure because the performance measure incorporates both precision and recall in its calculation. We use AUC because it has an a property or characterises of being insensitive to changes in the class distribution [20]. We apply 10 fold cross validation method for removing bias while evaluating the predictive models.

IV. EXPERIMENTAL DATASET

We conduct experiments on 226 publicly available web-services² shared on GitHub by the authors of the paper by Ouni et al. [6]. We conduct analysis on publicly available dataset so that our experimental results can be replicated and used for benchmarking and comparison. Juristo et al. and Kitchenham et al. work on experimental software engineering emphasizes the need of easy replication for increasing the reliability and validity of outcome in experiments [21][22]. Since our approach presented in this paper is the first work on applying several machine learning algorithms, feature ranking techniques and data sampling techniques, our objective is to

¹gromit.iiar.pwr.wroc.pl/p_inf/ckjm/

²<https://github.com/ouniali/WSantipatterns>

facilitate researchers to reproduce our work and use our results as baseline. The web-services used in our experiments belong to several domains such as education, media, weather, science, shipping, financial, search and travel. The web-service data is of high quality as the anti-patters are manually validated by the authors of the paper [6] who shared the dataset online on GitHub. The web-services belonging to the financial domain had the maximum number of anti-patterns.

Table I shows the number and percentages of 5 anti-patterns in our experimental dataset consisting of 226 web-services. We observe that the GOWS anti-pattern is present in 21 out of 226 web-services. There are 5.75% of web-services containing FGWS anti-pattern. The DWS, CWS and AWS anti-pattern are present in 14, 21 and 24 web-services respectively. As can be inferred from the Table I, the dataset has an imbalanced class distribution with respect to each anti-pattern. The dataset has an intrinsic imbalance characteristics due the nature of the dataspace. The number of web-services having a particular anti-pattern will naturally be in minority but not rare. Table I shows that an anti-pattern is minority (5% to 10%) but not rare (for example, lower than 1%). Also our dataset is small in-terms of the sample size and as shown in Table I, our problem is a machine learning problem consisting of a small sample size and imbalanced data [6][12] [11].

TABLE I: Number and Percentages of 5 Anti-Patterns in 226 Web-Services

	GOWS	FGWS	DWS	CWS	AWS
# AP	21	13	14	21	24
% AP	9.29	5.75	6.19	9.29	10.62

V. EXPERIMENTAL RESULTS AND FINDINGS

We follow a research methodology consisting of formulation of *four grounded research questions* guiding our scientific study. We apply the Sandberg et al. approach of formulating research questions based on existing work, literature and theories in our area [23]. Our main aim is to formulate research questions (RQ) and provide answers to them for facilitating development of tools and creation of methods for automatically detecting web-services anti-patterns through static analysis of source code.

RQ 1: Is there a statistically significant difference between the aggregate values of source code metrics computed at the web-service level between the two groups consisting of web-services containing a given anti-pattern and the web-services not containing the anti-pattern?

Our dataset consists of 226 web-services and our focus is on detecting five types of anti-patterns. For each of the five anti-patterns, we create two dataset. One dataset consists of web-services containing the anti-pattern and the other dataset consists of web-services which does not contain the anti-pattern. For example, for the GOWS anti-pattern, one dataset consists of 21 web-services and the other dataset consists of 105 web-services. For the FGWS anti-pattern, one dataset consists of 13 web-services and the other dataset consists of 213 web-services. In this experiment, the two datasets or the two samples are independent as there are no common web-services between the two datasets. We compute 148 values for each of the web-service consisting of aggregate and summary values like mean, median, standard deviation, Q1 and Q3 for each of the 18 object-oriented source code metrics. The number of values is 148 as there are some aggregate and

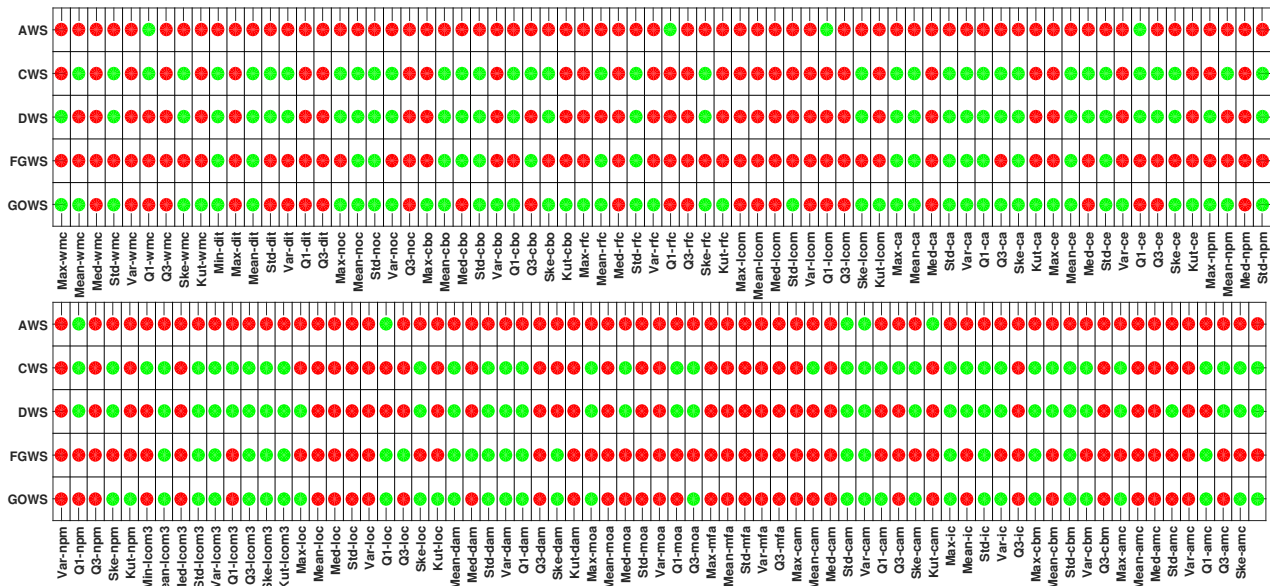


Fig. 2: Hypothesis Testing Results - based on p-value of source code metrics

summary values which cannot be computed (or have 0 or null values) for some of the metrics. We state the null hypothesis H_0 as a claim that there is no statistically significant difference in the mean value of the feature (independent variable) for the two datasets. The null hypothesis if accepted means that there is no influence of the metric or the independent variable on the dependent variable (dependent variable representing whether the web service contains the anti-pattern). The alternate hypothesis in this experiment H_a is the mathematical opposite of the null hypothesis H_0 which states that there is a statistically significant difference between the means of the metric value for the two dataset and hence the independent variable is potentially a predictor of the dependent variable.

We apply the Wilcoxon rank-sum test for hypothesis testing. Our objective is to determine if the null hypothesis should be accepted or rejected. We use the commonly used Wilcoxon rank-sum test methodology because the approach does not assume that the population follows a normal distribution. We use Wilcoxon rank-sum test as it is much more efficient than other types of parametric tests which assumes that the data (population and samples being analysed or compared) is normally distributed. Hence we do not make any parametric assumptions while conducting the null test. We set the significance level as 0.05 and compare the p-value with 0.05. We reject the null hypothesis (H_0) if the computed p-value is less than the cut-off value of 0.05. We reject the null hypothesis at a p-value of less than 0.05 as our inference is that the difference between the mean of the two groups is not because of chance and instead the difference in the mean is statistically significant. Figure 2 displays an illustration which reveals the cases for which the null hypothesis is accepted and rejected. A red dot denotes that the null hypothesis is rejected and a green dot denotes that the null hypothesis is accepted. Figure 2 reveals several metrics for which the number of red dots are ≥ 3 . Our analysis reveals that there are 91 metrics having red dots ≥ 3 . We set the cutoff (a heuristic) as number of red dots ≥ 3 and use 91 metrics as one set of features (Reduce Metric - RM) for the learning algorithm. We discard the remaining 57 metrics as non-informative or irrelevant based on the statistical significance test results. We conduct experiments with both the sets: AM (All 148 metrics) and RM (A reduced set of 91 metrics). Following is the answer to RQ1.

Answer RQ1: *Statistical significance test reveals that out of the 148 source code metrics as features, 91 metrics influences anti-pattern and there is a relation between 91 metrics and occurrence of the five types of anti-patterns.*

Feature Ranking Detailed Results : Table II displays the output of all the feature ranking technique. Table II displays the rank of every feature across all the anti-patterns for every feature ranking technique. We select (and highlight using shaded cell in Table II) top k attributes which is a logarithmic function of the total number of features. Table II reveals the ranking of each feature where the value in the cell represents the rank of the given feature. Table II shows that the rank

of the feature SKE_RFC is 1 for the GOWS and Information Gain Combination. We observe that there are some feature which are ranked highly across anti-patterns type and across feature selection techniques. We also notice that there are some attributes which are ranked higher for a particular anti-pattern or ranking technique. Table II shows that different feature selection techniques produces different output. An interesting finding is that few source code metrics have a higher rank for some anti-patterns while the same metrics is ranked lower for other anti-pattern within the context of a particular feature selection technique.

Detailed Results of 960 Predictive Models: Tables III, IV and V presents a detailed analysis and result of all the experimental executions for every combination of the experimental parameters. We are using 4 data sampling techniques, 8 learning algorithms 5 feature ranking techniques (total 6 after including the all metrics case) and 5 different types of anti-patterns. Thus the total number of unique classification models built by us are $4 \times 8 \times 6 \times 5 = 960$. We use 10 fold cross validation while evaluating the accuracy of each classifier. Table III shows that the accuracy of LB classifier for the case of Reduced Metrics, WTOSMP sampling technique and AWS anti-pattern is 88.05. Table III reveals a high f-measure value (0.94) for the NB and MLP classifier for the case of AWS anti-pattern and WTOSMP sampling technique. Table IV presents detailed results for the all-metrics case. Table IV shows that the f-measure value for NB and LOG classifier is 0.93 for the case of WTOSMP sampling technique and GOWS anti-pattern. From Table IV, we infer that the accuracy of RF classifier is above 97% for SMOTE and FGWS and DWS anti-pattern. Table V provided detailed result across all feature ranking methods. From Table V, we infer that the f-measure value of the RF classifier is above 0.9 for several combinations of data sampling technique and feature ranking techniques.

Distribution Comparison and Visualization using Boxplot: Tables VI, VII and VIII shows the descriptive statistics of the performance of the various combinations of data sampling, feature ranking and learning algorithms. Figures 3, 4 and 5 shows the boxplots for comparing the degree of dispersion, interquartile range, outliers and skewness in the accuracy, f-measure and AUC values for all the combinations of techniques. The red line in the boxplot of Figures 3, 4 and 5 marks the median point of the data (mid-point) and the red line thus divides the box into two segments. Figure 3 shows that the median AUC value of NB and RF is higher than the corresponding values for other learning algorithms such as LOG, J48, MLP, LB, ADA and BAG. Figure 3 shows that the median value for LOG is the lowest in comparison to other classification algorithms. From the distributional characteristics of data shown the bog-plots of Figure 4, we infer that the inter-quartile range representing the middle box signifying the middle 50% of the values of the given variable range for LOG and BAG is more than the corresponding

TABLE IV: Performance Parameters (All Metrics)

		Accuracy								F-Measure							
		J48	RF	LB	ADA	BAG	MLP	NB	LOG	J48	RF	LB	ADA	BAG	MLP	NB	LOG
WTOSMP	GOWS	92.04	92.04	92.92	92.48	93.36	92.48	88.05	87.61	0.96	0.96	0.96	0.96	0.96	0.96	0.93	0.93
WTOSMP	FGWS	93.81	95.58	95.13	94.69	94.25	94.69	61.06	88.5	0.97	0.98	0.97	0.97	0.97	0.97	0.74	0.94
WTOSMP	DWS	95.13	97.35	96.46	96.02	94.69	95.58	93.36	93.36	0.97	0.99	0.98	0.98	0.97	0.98	0.98	0.96
WTOSMP	CWS	91.59	92.48	92.48	92.92	92.04	92.04	81.42	85.84	0.95	0.96	0.96	0.96	0.96	0.96	0.89	0.92
WTOSMP	AWS	87.17	88.05	83.63	89.38	89.38	86.28	57.08	79.65	0.93	0.94	0.91	0.94	0.94	0.92	0.71	0.88
RANSMP	GOWS	80	83.67	81	80.33	84.33	82.67	82.67	74.67	0.85	0.88	0.86	0.85	0.88	0.87	0.86	0.81
RANSMP	FGWS	77.3	81.08	77.3	72.97	72.97	71.35	71.35	67.03	0.81	0.86	0.82	0.78	0.8	0.76	0.75	0.71
RANSMP	DWS	92.5	93	93.5	93.5	92	91	92	87.5	0.94	0.94	0.95	0.95	0.94	0.93	0.94	0.9
RANSMP	CWS	79.33	85	78.67	78	79.67	79.33	80	75.33	0.84	0.88	0.83	0.83	0.84	0.84	0.83	0.81
RANSMP	AWS	56.52	66.09	60	59.71	60.58	60.58	53.91	60.29	0.65	0.76	0.7	0.71	0.72	0.69	0.63	0.69
DOWNSMP	GOWS	77.04	82.47	80.71	78.83	81.05	80.82	83.82	69.46	0.8	0.85	0.83	0.82	0.83	0.83	0.85	0.74
DOWNSMP	FGWS	74.46	81.84	77.68	75.72	75.28	73.98	71.27	70.69	0.8	0.87	0.83	0.82	0.82	0.8	0.75	0.76
DOWNSMP	DWS	91.2	90.56	91.14	90.67	91.48	88.59	92.09	83.62	0.93	0.93	0.93	0.93	0.94	0.91	0.94	0.87
DOWNSMP	CWS	77.48	84.29	82.24	82.94	78.49	80.25	84.67	73.82	0.8	0.86	0.85	0.85	0.81	0.82	0.86	0.76
DOWNSMP	AWS	61.69	57.67	61.83	61.44	54.67	60.78	47.44	54.26	0.64	0.6	0.64	0.64	0.53	0.63	0.48	0.55
SMOTE	GOWS	89.63	93.7	89.63	90.74	91.48	91.85	91.11	91.48	0.93	0.95	0.93	0.93	0.94	0.94	0.93	0.94
SMOTE	FGWS	94.35	97.58	96.37	89.52	96.37	95.16	69.35	87.5	0.96	0.98	0.98	0.94	0.98	0.97	0.77	0.92
SMOTE	DWS	95.33	97.67	96.89	96.89	96.5	97.28	97.67	97.28	0.97	0.99	0.98	0.98	0.98	0.98	0.98	0.98
SMOTE	CWS	93.49	96.93	94.64	93.1	92.34	95.4	88.51	94.64	0.95	0.98	0.96	0.95	0.94	0.97	0.91	0.96
SMOTE	AWS	81.96	88.63	84.71	78.04	84.31	90.98	60.78	85.49	0.86	0.91	0.88	0.83	0.88	0.93	0.63	0.88

TABLE V: Performance Parameters (Feature Ranking Methods)

		Accuracy								F-Measure							
		J48	RF	LB	ADA	BAG	MLP	NB	LOG	J48	RF	LB	ADA	BAG	MLP	NB	LOG
WTOSMP	FR1	92.48	93.81	92.92	92.48	93.36	92.48	91.15	91.59	0.96	0.97	0.96	0.96	0.96	0.96	0.95	0.95
WTOSMP	FR2	88.05	88.05	88.05	90.27	89.82	89.38	75.22	89.38	0.94	0.94	0.94	0.95	0.95	0.94	0.85	0.94
WTOSMP	FR3	91.59	92.48	94.69	93.36	90.27	93.81	89.38	92.48	0.96	0.96	0.97	0.96	0.95	0.97	0.94	0.96
WTOSMP	FR4	90.71	87.17	91.15	90.27	90.71	91.59	87.17	91.15	0.95	0.93	0.95	0.95	0.95	0.96	0.93	0.95
WTOSMP	FR5	91.59	92.48	91.59	92.04	92.92	92.48	89.82	92.48	0.95	0.96	0.95	0.96	0.96	0.96	0.94	0.96
RANSMP	FR1	78.7	81.99	81.84	81.74	83.99	83.08	84.92	81.33	0.82	0.85	0.84	0.84	0.86	0.86	0.87	0.84
RANSMP	FR2	72.8	70.55	71.18	73.6	74.06	78.09	77.38	77.65	0.75	0.73	0.74	0.76	0.77	0.8	0.78	0.8
RANSMP	FR3	80.16	80.74	79.33	78.95	79.45	82.33	83.06	81.8	0.82	0.83	0.82	0.82	0.82	0.85	0.86	0.84
RANSMP	FR4	74.21	66.82	69.3	71.53	72.09	75.63	70.29	74.15	0.75	0.71	0.73	0.76	0.76	0.77	0.77	0.78
RANSMP	FR5	77.8	80.77	79.12	80.58	82.04	81.27	86.93	80.03	0.81	0.83	0.82	0.83	0.84	0.84	0.89	0.82
DOWNSMP	FR1	88.45	93.14	91.7	90.98	89.17	91.7	87.73	91.34	0.92	0.95	0.94	0.94	0.92	0.94	0.91	0.94
DOWNSMP	FR2	85.15	89.49	85.15	84.78	85.87	88.41	77.9	82.97	0.89	0.92	0.9	0.9	0.9	0.92	0.82	0.88
DOWNSMP	FR3	87.46	89.25	89.25	87.81	89.61	92.83	89.25	91.76	0.91	0.92	0.93	0.91	0.93	0.95	0.93	0.94
DOWNSMP	FR4	87.19	88.61	85.41	80.78	83.99	85.05	78.29	86.12	0.91	0.92	0.9	0.87	0.89	0.89	0.85	0.9
DOWNSMP	FR5	89.82	90.55	90.55	88.73	91.27	90.91	88.73	92	0.93	0.93	0.93	0.92	0.94	0.94	0.92	0.94
SMOTE	FR1	82.53	83.1	81.13	83.45	83.74	81.31	86.32	80.14	0.85	0.86	0.84	0.86	0.86	0.84	0.88	0.83
SMOTE	FR2	73.61	73.38	71.89	75.39	74.9	80.25	78.11	77.5	0.77	0.77	0.75	0.78	0.78	0.82	0.78	0.81
SMOTE	FR3	75.34	78.17	75.58	79	78.69	81.41	82.71	80.57	0.79	0.81	0.79	0.82	0.82	0.84	0.86	0.84
SMOTE	FR4	73.66	70.47	72.65	74.21	75.05	74.24	72.42	76.29	0.76	0.75	0.76	0.78	0.78	0.77	0.79	0.8
SMOTE	FR5	78.68	82.92	81.42	80.86	83.99	82.82	86.07	83.17	0.81	0.85	0.83	0.83	0.86	0.85	0.88	0.85

values for NB, J48, RF, MLP, LB and ADA. Similarly, we observe that the inter-quartile range for RANSMP is more than the inter-quartile ranges of WTOSMP, SMOTE and DOWNSMP. Box-plots in Figure 3 and 4 are for all metrics and reduced metrics whereas the box-plot in Figure 5 is overall and for all the combinations. From Figure 5, we observe that the boxplot for the AUC values for LOG, LB, SMOTE and FR1 are comparatively taller. This means that there is more variation in the AUC values obtained from multiple executions. While the box-plots in Figures 3, 4 and 5 shows the results in the form of a visualization, Tables VI, VII and VIII presents exact values obtained from the experiments for all the combinations of the experimental parameters.

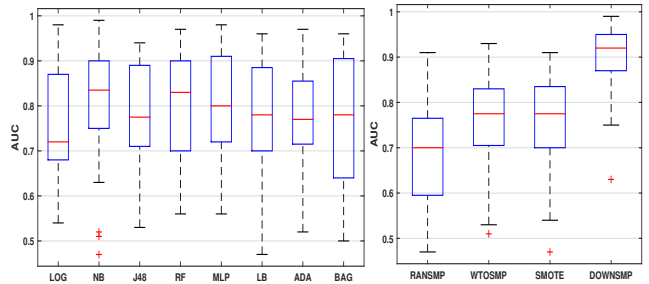


Fig. 3: Boxplot: AUC (All Metrics)

TABLE VI: Descriptive Statistics of the Performance in-terms of AUC (All Metrics)

AUC							
	Min	Max	Mean	Median	Std Dev	Q1	Q3
LOG	0.54	0.98	0.76	0.72	0.13	0.68	0.87
NB	0.47	0.99	0.79	0.84	0.15	0.75	0.90
J48	0.53	0.94	0.77	0.78	0.12	0.71	0.89
RF	0.56	0.97	0.80	0.83	0.13	0.70	0.90
MLP	0.56	0.98	0.79	0.80	0.13	0.72	0.91
LB	0.47	0.96	0.77	0.78	0.13	0.70	0.89
ADA	0.52	0.97	0.76	0.77	0.13	0.72	0.86
BAG	0.50	0.96	0.75	0.78	0.15	0.64	0.91
RANSMP	0.47	0.91	0.69	0.70	0.11	0.60	0.77
WTOSMP	0.51	0.93	0.75	0.78	0.12	0.71	0.83
SMOTE	0.47	0.91	0.75	0.78	0.11	0.70	0.84
DownSMP	0.63	0.99	0.90	0.92	0.07	0.87	0.95

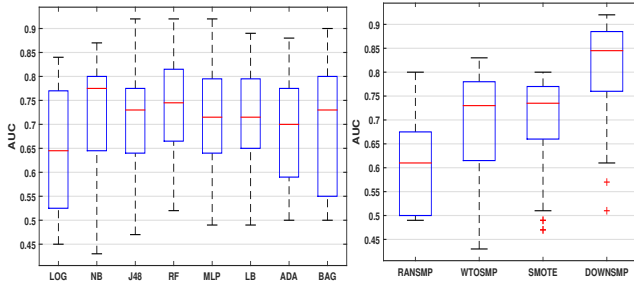


Fig. 4: Boxplot: AUC (Reduced Metrics)

TABLE VII: Descriptive Statistics of the Performance in-terms of AUC (Reduced Metrics)

AUC							
	Min	Max	Mean	Median	Std Dev	Q1	Q3
LOG	0.45	0.84	0.66	0.65	0.13	0.53	0.77
NB	0.43	0.87	0.71	0.78	0.13	0.65	0.80
J48	0.47	0.92	0.71	0.73	0.13	0.64	0.78
RF	0.52	0.92	0.74	0.75	0.12	0.67	0.82
MLP	0.49	0.92	0.71	0.72	0.13	0.64	0.80
LB	0.49	0.89	0.71	0.72	0.11	0.65	0.80
ADA	0.50	0.88	0.68	0.70	0.12	0.59	0.78
BAG	0.50	0.90	0.69	0.73	0.15	0.55	0.80
RANSMP	0.49	0.80	0.61	0.61	0.09	0.50	0.68
WTOSMP	0.43	0.83	0.69	0.73	0.11	0.62	0.78
SMOTE	0.47	0.80	0.69	0.74	0.10	0.66	0.77
DownSMP	0.51	0.92	0.81	0.85	0.11	0.76	0.89

Null Hypothesis Statistical Significance Testing: We apply the nonparametric Wilcoxon signed-rank test for statistical hypothesis testing (with a Bonferroni correction) to investigate if the two samples consisting of AUC values generated from two different combinations of experimental parameters were selected from population having same or different distributions. The nonparametric Wilcoxon signed-rank test is a type of paired t-test used to check if the null hypothesis can be accepted or rejected. Figures 6, 7 and 8 shows the output of the hypothesis testing consisting of a comparison between four data sampling methods, eight classification algorithms and six feature ranking methods on AUC performance metrics. Figures 6, 7 and 8 consists of a red dot or a green dot. A red dot means that H_0 (null hypothesis) is rejected and a green dot means

TABLE VIII: Descriptive Statistics of the Performance in-terms of AUC (Overall)

AUC							
	Min	Max	Mean	Median	Std Dev	Q1	Q3
LOG	0.44	0.94	0.67	0.68	0.14	0.52	0.79
NB	0.46	0.95	0.73	0.76	0.13	0.69	0.82
J48	0.49	0.96	0.72	0.75	0.13	0.60	0.82
RF	0.54	0.98	0.76	0.76	0.11	0.68	0.86
MLP	0.49	0.98	0.71	0.73	0.13	0.59	0.81
LB	0.50	0.99	0.73	0.74	0.12	0.65	0.81
ADA	0.50	0.98	0.71	0.73	0.12	0.60	0.80
BAG	0.50	0.98	0.72	0.75	0.14	0.60	0.83
RANSMP	0.44	0.90	0.70	0.71	0.12	0.59	0.80
WTOSMP	0.44	0.90	0.72	0.75	0.12	0.61	0.81
SMOTE	0.44	0.92	0.65	0.64	0.12	0.54	0.76
DownSMP	0.44	0.90	0.68	0.71	0.14	0.54	0.80
FR1	0.44	0.99	0.74	0.76	0.13	0.65	0.83
FR2	0.44	0.98	0.70	0.71	0.13	0.59	0.81
FR3	0.47	0.96	0.72	0.76	0.13	0.62	0.82
FR4	0.46	0.97	0.71	0.73	0.13	0.60	0.81
FR5	0.47	0.96	0.72	0.74	0.12	0.63	0.81

TABLE IX: Mean Difference AUC (Data Sampling Methods)

	WTOSMP	RANSMP	DownSMP	SMOTE
WTOSMP	0.000	-0.086	-0.084	-0.193
RANSMP	0.086	0.000	0.002	-0.107
DownSMP	0.084	-0.002	0.000	-0.109
SMOTE	0.193	0.107	0.109	0.000

that H_0 is accepted. H_0 being accepted means that there is no statistically significant difference between the two techniques or combination of parameters under comparison. While performing statistical significance testing, we use a standard cut-off value of 0.05 (signifying that the null hypothesis is rejected when $p \leq 0.05$ and not rejected when $p > 0.05$).

According to Figure 6, there are a total of 6 possible comparisons as there are 4 data sampling methods. We observe from Figure 6 that 5 out of 6 cells contains red dots and only 1 cell contains a green dot. Figure 6 reveals that there is no statistically significant difference between DownSMP and RANSMP. However, there is a statistically significant difference between the performance of SMOTE and WTOSMP, SMOTE and RANSMP, SMOTE and DownSMP, DownSMP and WTOSMP, RANSMP and WTOSMP. Figure 7 shows the pairwise comparison between the eight classification techniques. Figure 7 reveals that there is a statistical significant difference between LOG and the remaining 7 classifiers. We observe that NB shows difference with LOG and not any other classification algorithm. We notice that RF and J48 performance are different based on the statistical hypothesis testing results as shown by the red dot at their intersecting cell in the Figure 7. Figure 7 shows that classification algorithm influences the accuracy results. Figure 8 displays the pairwise comparison results for the 6 feature ranking methods. Figure 7 reveals there is a statistical significant difference between the performances between several feature ranking methods as shown by the red dots in the intersecting cell. From Figure 7, we infer that the feature ranking method FR1 results in statistically significant

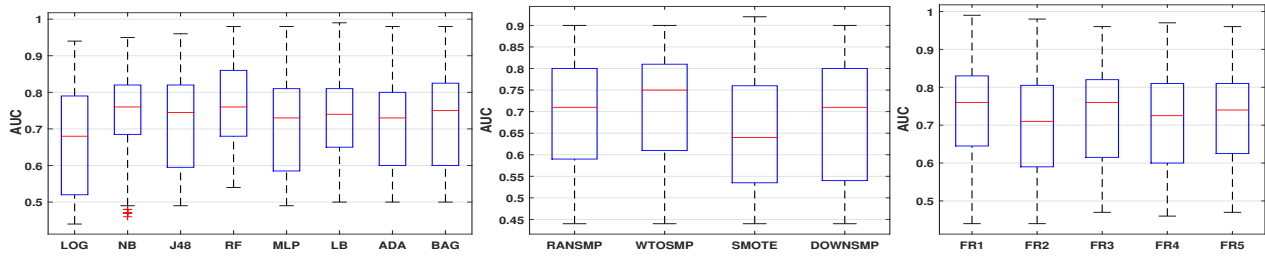


Fig. 5: Boxplot: AUC (Overall)

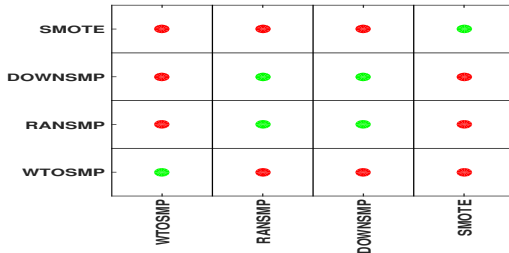


Fig. 6: Wilcoxon Signed Rank Test Analysis: P-Value of AUC (Data Sampling Methods)

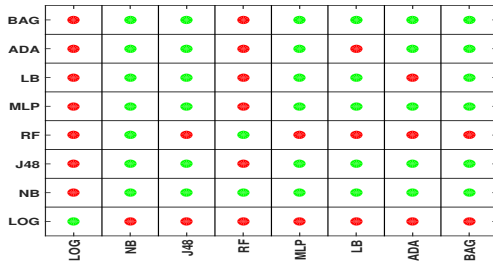


Fig. 7: Wilcoxon Signed Rank Test Analysis: P-Value of AUC (Classification Techniques)

different performance than the other feature ranking methods. We do not observe any difference in performance between FR4 and FR2 as well as FR4 and FR3.

Mean AUC Value Difference: Tables IX, X and XI shows the difference between the mean AUC values for various data sampling techniques, feature ranking techniques and classification algorithms. A higher AUC value means a better performance. The cells in Tables IX, X and XI displays the difference between the mean AUC value of the technique mentioned in the row heading and the technique mentioned in the column heading. A positive value of the cell indicates that the technique mentioned in the row heading outperforms the technique mentioned in the column heading. From Table IX, we infer that the NB outperforms LOG classification algorithm in-terms of the AUC value by 0.056. Table IX shows that NB outperforms all other classification techniques except RF. Table IX reveals that RF performs the best as its mean AUC value is higher than the mean AUC values

TABLE X: Mean Difference AUC (Classification Algorithm)

	LOG	NB	J48	RF	MLP	LB	ADA	BAG
LOG	0.000	-0.056	-0.043	-0.086	-0.044	-0.054	-0.031	-0.039
NB	0.056	0.000	0.013	-0.030	0.013	0.003	0.025	0.017
J48	0.043	-0.013	0.000	-0.043	0.000	-0.010	0.012	0.004
RF	0.086	0.030	0.043	0.000	0.042	0.032	0.055	0.047
MLP	0.044	-0.013	0.000	-0.042	0.000	-0.010	0.012	0.005
LB	0.054	-0.003	0.010	-0.032	0.010	0.000	0.022	0.015
ADA	0.031	-0.025	-0.012	-0.055	-0.012	-0.022	0.000	-0.008
BAG	0.039	-0.017	-0.004	-0.047	-0.005	-0.015	0.008	0.000

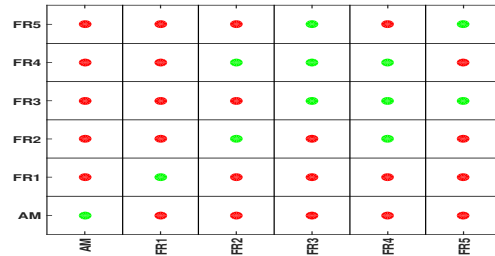


Fig. 8: Wilcoxon Signed Rank Test Analysis: P-Value of AUC (Feature Ranking Methods)

of all other classification algorithms. We observe that the mean AUC value of LOG is lower than the mean AUC values of all other classification algorithms. The highest difference between the mean AUC values of classification algorithms is between RF and LOG which is 0.086. Table X shows that SMOTE is the best performing data sampling method in comparison to other sampling methods used in our experiments. The difference in the mean AUC values of SMOTE and WTOSMP is the highest which is 0.193. There is a minor difference between the overall performance of DOWNSMP and RANSMP. The second best performing technique out of the four techniques is RANSMP. Table XI reveals that AM results in the best performance in comparison to the five feature ranking methods. All the values in the cells for which AM is in the row heading are positive. The maximum difference is between AM and FR4. We observe that FR2 is the weakest approach in comparison to all other approaches. There is a minor difference in the performance of FR2 and FR4.

RQ 2: Is there a statistically significant difference in the prediction performance of five different feature selection techniques and four different data sampling techniques?

TABLE XI: Mean Difference AUC (Feature Ranking Methods)

	AM	FR1	FR2	FR3	FR4	FR5
AM	0.000	0.038	0.072	0.053	0.067	0.049
FR1	-0.038	0.000	0.034	0.016	0.029	0.012
FR2	-0.072	-0.034	0.000	-0.019	-0.005	-0.023
FR3	-0.053	-0.016	0.019	0.000	0.013	-0.004
FR4	-0.067	-0.029	0.005	-0.013	0.000	-0.017
FR5	-0.049	-0.012	0.023	0.004	0.017	0.000

Which of the feature selection and data sampling technique(s) yields the best result?

Answer RQ 2: *There is a statistically significant difference in the performance of various feature selection techniques and data sampling techniques. SMOTE is the best performing data sampling method in comparison to other sampling methods used in our experiments. OneR feature ranking outperforms all other feature ranking techniques.*

RQ 3: What is the relative performance of the eight different types of learning algorithms and classifiers in-terms of the AUC, accuracy and f-measure metrics? Is there a statistically significant difference in the prediction performance of eight different learning algorithms or classifiers?

Answer RQ 3: *There is a statistically significant difference in the performance of the eight different types of learning algorithms. Random Forest learning algorithms performs best followed by Nave Bayes for the task of web-service anti-pattern detection using source code metrics.*

VI. CONCLUSION

Our main conclusion is that it is possible to predict the occurrence of web-service anti-patterns using object oriented source code metrics derived from the source code of the web services and machine learning classifiers. We observe a correlation between the value of certain source code metrics and web-services anti-patterns. We conduct a series of experiments by building 960 classification models. Our conclusion is that data sampling method is an important pre-processing step since the class distribution is imbalanced and there is a statistically significant difference between the performances of various data sampling methods. Our results indicate that SMOTE is the best performing data sampling method. Since our dataset consists of several dimensions, we observe that feature ranking method is useful to reduce the dimension of the dataset and identify informative features. Our results reveal that OneR feature ranking method outperforms all other feature ranking method used in our experiments. In terms of the learning algorithms, we conclude that there is a statistically significant difference between the performance of the 8 classifiers and Random Forest approach outperforms all other classification algorithms.

REFERENCES

- [1] Y. Kun, W. Xiao-Ling, and Z. Ao-Ying, "Underlying techniques for web services: A survey," in *Journal of software*. Citeseer, 2004.
- [2] E. Newcomer and G. Lomow, *Understanding SOA with Web services*. Addison-Wesley, 2005.
- [3] F. Curbera, M. Duftler, R. Khalaf, W. Nagy, N. Mukhi, and S. Weerawarana, "Unraveling the web services web: an introduction to soap, wsdl, and uddi," *IEEE Internet computing*, vol. 6, no. 2, pp. 86–93, 2002.
- [4] J. Kral and M. Zemlicka, "The most important service-oriented antipatterns," in *Software Engineering Advances, 2007. ICSEA 2007. International Conference on*. IEEE, 2007, pp. 29–29.
- [5] J. Král and M. Žemlicka, "Popular soa antipatterns," in *Future Computing, Service Computation, Cognitive, Adaptive, Content, Patterns, 2009. COMPUTATIONWORLD'09. Computation World.*. IEEE, 2009, pp. 271–276.
- [6] A. Ouni, M. Kessentini, K. Inoue, and M. O. Cinnéide, "Search-based web service antipatterns detection," *IEEE Transactions on Services Computing*, 2015.
- [7] A. Ouni, R. Gaikovina Kula, M. Kessentini, and K. Inoue, "Web service antipatterns detection using genetic programming," in *Proceedings of the 2015 Annual Conference on Genetic and Evolutionary Computation*. ACM, 2015, pp. 1351–1358.
- [8] F. Palma, N. Moha, G. Tremblay, and Y.-G. Guéhéneuc, "Specification and detection of soa antipatterns in web services," in *European Conference on Software Architecture*. Springer, 2014, pp. 58–73.
- [9] F. Palma, M. Nayrolles, N. Moha, Y.-G. Guéhéneuc, B. Baudry, and J.-M. Jézéquel, "Soa antipatterns: An approach for their specification and detection," *International Journal of Cooperative Information Systems*, vol. 22, no. 04, p. 1341004, 2013.
- [10] N. Moha, F. Palma, M. Nayrolles, B. J. Conseil, Y.-G. Guéhéneuc, B. Baudry, and J.-M. Jézéquel, "Specification and detection of soa antipatterns," in *ICSOC*. Springer, 2012, pp. 1–16.
- [11] C. Mateos, M. Crasso, A. Zunino, and J. L. O. Coscia, "Detecting wsdl bad practices in code-first web services," *International Journal of Web and Grid Services*, vol. 7, no. 4, pp. 357–387, 2011.
- [12] J. L. Ordiales Coscia, C. Mateos, M. Crasso, and A. Zunino, "Anti-pattern free code-first web services for state-of-the-art java wsdl generation tools," *International Journal of Web and Grid Services*, vol. 9, no. 2, pp. 107–126, 2013.
- [13] M. Nayrolles, F. Palma, N. Moha, and Y.-G. Guéhéneuc, "Soda: A tool support for the detection of soa antipatterns," in *International Conference on Service-Oriented Computing*. Springer, 2012, pp. 451–455.
- [14] S. R. Chidamber and C. F. Kemerer, "A metrics suite for object oriented design," *IEEE Transactions on software engineering*, vol. 20, no. 6, pp. 476–493, 1994.
- [15] S. Chidamber and C. F. Kemerer, *Towards a metrics suite for object oriented design*. ACM, 1991, vol. 26, no. 11.
- [16] 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.
- [17] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Transactions on knowledge and data engineering*, vol. 21, no. 9, pp. 1263–1284, 2009.
- [18] H. Liu and L. Yu, "Toward integrating feature selection algorithms for classification and clustering," *IEEE Transactions on knowledge and data engineering*, vol. 17, no. 4, pp. 491–502, 2005.
- [19] T.-S. Lim, W.-Y. Loh, and Y.-S. Shih, "A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms," *Machine learning*, vol. 40, no. 3, pp. 203–228, 2000.
- [20] M. Sokolova and G. Lapalme, "Performance measures in classification of human communications," *Advances in Artificial Intelligence*, pp. 159–170, 2007.
- [21] N. Juristo and O. S. Gómez, "Replication of software engineering experiments," in *Empirical software engineering and verification*. Springer, 2012, pp. 60–88.
- [22] B. Kitchenham, "The role of replications in empirical software engineering word of warning," *Empirical Software Engineering*, vol. 13, no. 2, pp. 219–221, 2008.
- [23] J. Sandberg and M. Alvesson, "Ways of constructing research questions: gap-spotting or problematization?" *Organization*, vol. 18, no. 1, pp. 23–44, 2011.