

A Survey on Class Imbalance Learning Algorithms

R. BuliBabu*, Dr. Mohammed Ali Hussain**

*Research Scholar, Department of Computer Science, Bharathiar University, Coimbatore, India.

**Professor, Dept. of Electronics and Computer Engineering, KLEF University, India

Email: rsmbabu@yahoo.com, alihussain.phd@gmail.com

ABSTRACT: We have gone through a systematic and comparative analysis of using learning methods on imbalanced data, by using different 11 learning algorithms on real word datasets from different types of applications and applicational domains. Our objective is to give practical approach to machine learning research to build different type of classifiers on class imbalanced datasets, and to guide the researchers some various possible guidelines for work. Our works proposes to analysis class imbalanced from a wide scope, difference learners, performing sampling methods, and applying performance measuring on various datasets.



INTRODUCTION

In the native world various domain classification occurs majorly from one of the various classes. In binary classification, the research typically chooses marginal positive class of his own interest. The imbalance in class distribution causes learning machine algorithms to act very poorly on the marginal class. In calculation, the cost of misclassifying on marginal class is usually much upper than the cost in various other misclassifications. Therefore a common query occurs for machine learning researcher to identify and improve the performance of classifiers when certain classes are relatively rare.

The answer is to sample the data, either random or knowledge based, for obtaining an alternative distribution class. Various methods and techniques are proposed, it was uncertain which technique gives best results on work. Number of researchers experimentally solved with the use of sampling using learning from class imbalanced datasets.

The tested evaluated in this study gives a completed background work, which includes two or more learners and datasets. In our work, we have taken 36 various datasets (section 1.1), Sampling techniques of 7, various learning algorithms 11. The proposed works in Section 1.5, a total of 1,232,000 classifiers were built in our experiments. In added, we also performed SPSS analysis for gaining variance (ANOVA) for understanding the statistical significance of the results. The variance analysis made our work very comprehensive and dynamically increased the reliability of our work.

We powerfully believe reliability and statistical validation give the strength and weakness of various methods and algorithms on real world domain applications.

BACKGROUND

Matwin et al, suggested a technique named one sided selection. One sided selection applied to knowledge based under sample to the best class by eliminating majority class, to eliminate either redundant or noisy. For the improve of performance on random resampling.

Barandela et al. proposed a **KNN** technique with $k = \text{integer value}$ which removes all main class that are misclassified using classifier. He also suggested to re-write calculation distance which case identify positives from -ve ones.

Chawla et al. narrated an oversampling intelligent method called Synthetic Minority Oversampling Technique (SMOTE). SM adds new and artificial minority by exploring between re-existing minority objects rather than eliminating duplicates from original examples. This technique will find the K-Nearest Neighbors first from minority class for every minority (it is recommended to take $k=5$). The sample examples then evaluate some or all NN, varying on the type of oversampling.

Han et al. gives a modified SMOTE method which is called borderline-smote . BSM which identifies minority examples, it considered the border of various decision regions of minority class in the space and then evaluate for SMOTE for oversampling the instance, instead of oversampling them using random data subset.

Jo & Japkowicz, suggested the implementation of cluster based over-sampling for class within balance and imbalance classes. Subsets of a class isolates the feature space of balanced class, and subsets of within class imbalance are called disjoints small. Small disjoints causes or laydown classifier performance, so our aim is to remove such data

We have performed RUS for various rates of 5%, 10%, 25%, 50%, 75%, and 90% for common classes, where SM, BSM and ROS have gone through oversampling performance at the rate of 50%, 100%,

200%, 300%, 500%, 750%, and 1000%. The performance has been done on Euclidean weights for both weighted and un-weighted. The total of combinations of 30 sampling techniques and parameters were utilized. .

PROPOSED WORK

In this part we show a brief explanation for eleven classification algorithms with various examples using different parameters in our practical approach. These classification algorithms are used basically in machine learning for class imbalanced dataset.

The learning methods are developed in WEKA. These methods are provided with default experimental values for checking the performance on the classifiers for all datasets based on analysis.

The learners which were used in the paper was built in a tool WEKA, by changing the mandatory value parameter which were done on experimental bases, which also show the advanced improvement in the performance of the classifiers of all datasets used for basic analysis

We have used dual version of decision tree learner namely C4.5D and C4.5N, these methods were developed using J48 in WEKA. C4.5D uses WEKA parameter default values, on the other C4.5N has provided with smoothing and pruning dataset values.

K-NN classifiers methods have been developed in WEKA, uses K=3 and K=5, which is denoted 3 or 4NN and 5NN. The parameters are set with the weighting distance which is mentioned as weight by inverse of distance for determining the classifying instance. In Naïve Bayes (NB) classifier, the default parameters were assumed from left.

Multilayer perceptions learners uses two different parameters which can change the training process. They are hidden layers and validationSetSize. The first one the hidden layer which contains three

nodes and it can be changed to 3 and the second one the ValidationSetSize has to be changed the 10% training data set to determine and signify when the process to stop the iteration

Artificial Neural networks uses a network called Radial basis function networks and it set 'numClusters' to 10.

Another rule based learner in Artificial Neural network was Ripper which uses the default parameters in each and every experiment and no changes were made to the default parameters which is denoted by Logistic regression LR

To construct decision trees ,the classifier The Random forest (RF) uses Random subspace method which produces the atmost prediction and there is no change in the default parameters and it should also uses the Support vector machine (SMO) and it is recongnised as SVM in this experiment which is linear kernel.Two parameters were used in this experiment,the complexity control 'c' changes from 1.0 to 5.0 and the second one 'build logistic models' which is set to be true

EXPERIMENTAL DATASETS

Table 1 stated the 35 datasets in our experimental study. The percentage varies from 1.33%(highly imbalanced) to 35% (only slightly imbalanced).Among the various application domains of the datasets 19 taken from UCI repository. Dr.Nitish Chawla provided Mammography datasets.Among them 15 datasets taken from the domain of SE measurements. Totally 214 examples (Glass-3) had smallest datasets, largest datasets each contain 20,000 observations. A binary class which is used to transform all the datasets and the binary classification problem is used in this experiment

SAMPLING TECHNIQUES

This part gives a brief overview of 7 sampling techniques

Table 1. Surveyed Datasets

Name	# minority	% minority	#attr
Sp3	47	1.33%	43
Sp4	92	2.31%	43
Mammography	260	2.32%	7
Nursery-3	328	2.53%	9
Solar flare -f	51	3.67	13
Letter -a	789	3.95	17
Car-3	69	3.99	7
Sp2	189	4.75	43
Cccs 12	16	5.67	9
SP1	229	6.28	43
pc1	76	6.87	16
Mw1	31	7.69	16
Glass-3	17	7.94	10
Kc3	43	9.39	16
Cm1	48	9.50	16
Cccs-8	27	9.57	9
Pendlqlts-5	1055	9.60	17
Satlmaqe-4	626	9.73	37
Optdlqlts-8	554	9.86	65
e-coll-4	35	10.52	8
Segment-5	330	14.29	20
Kc1	325	15.42	16
Lm1	1687	19.06	16
Letter vowel	3878	19.39	17
Cccs-4	55	19.60	9
Kc2	106	20.38	16

Contra-2	333	22.69	10
Splcejunk-2	768	24.08	61
Vehicle-1	212	25.06	19
Haberman	81	26.47	4
Yeast-2	429	28.91	9
Phoneme	1525	29.35	6
Ccs2	83	29.43	9
German credit	300	30.00	21
Plma diabetes	268	34.90	9

Totally seven sampling techniques were used in this section 1. Random oversampling(ROS) 2. One sided selection (OSS) 3. Wilson’s editing(WE) 4. Random undersampling(RUS) 5. SMOTE (SM)6. Borderline SMOTE (BSM). These requires a value which is to be set to a parameter. Example ROS300 means random oversampling with the parameter 300(all the sampling techniques were explained in the below table

Random Minority oversampling (ROS) and Random majority under sampling(RUS) are the two important preprocessing techniques. Majority class of instances are discarded from the dataset and the minority class are randomly duplicated

EXPERIMENTAL DESIGN

Summary can be taken from our experiment and it can be discussed as follows. Among 35 datasets , 20 five –fold cross validation(CV) were executed. To run CV , each iteration consists of four folds and the one fold maintains to test dataset. Totally 31 sampling techniques were applied to the data and 11 different learners was finalized on the dataset which was tested on dataset(depend on CV)

Table 2. SVM $\pi < 10\%$ - Sampling

Level	AUC	HSD	Level	G	HSD
ROS1000	0.898	A	RUS5	82.24	A
RUS5	0.897	AB	CBOS	80.36	AB
SM1000	0.890	AB	ROS1000	76.49	CB
BSM1000	0.886	AB	SM1000	75.17	CB
CBOS	0.872	B	BSM1000	71.93	C
WE-W	0.821	C	OSS	51.81	D
OSS	0.818	C	WE-W	45.28	E
NONE	0.809	C	NONE	41.75	E
Level	AUC	HSD	Level	G	HSD
ROS1000	0.861	A	ROS10001	78.156	A
SM300	0.860	A	SM1000	78.017	A
BSM300	0.856	A	RUS5	76.431	AB
RUS25	0.849	AB	BSM1000	75.851	AB
CBOS	0.830	CB	CBOS	73.173	B
WE-W	0.828	C	WE-W	51.725	C
OSS	0.818	CD	OSS	45.977	D
NONE	0.798	D	NONE	42.505	D

Table 3. RF -Sampling, $\pi < 10\%$

Level	AUC	HSD	Level	G	HSD
RUS5	0.892	A	RUS5	83.08	A
SM1000	0.865	B	SM1000	64.16	B
BSM1000	0.859	BC	BSM1000	60.17	BC
ROS-300	0.847	BCD	ROS1000	59.47	BC
WE-W	0.842	BCD	CBOS	57.20	CD
NONE	0.837	CD	OSS	56.90	CD
CBOS	0.825	DE	WE-E	51.69	DE
OSS	0.810	E	NONE	49.08	E
Level	AUC	HSD	Level	G	HSD
RUS10	0.862	A	RUS10	79.16	A
SM750	0.857	AB	SM1000	70.54	B
BSM1000	0.852	AB	BSM1000	68.80	BC
WE-W	0.846	AB	WE-W	65.29	CD
ROS200	0.844	ABC	ROS300	64.33	DE
OSS	0.839	BC	OSS	61.65	DEF
NONE	0.836	BC	CBOS	61.24	EF
CBOS	0.825	C	NONE	59.76	F

CV runs on datasets of overall 20 fivefold, 3500 different training datasets. 31 sampling techniques, and no sampling, applied on training datasets, resulting in $32 \times 3500 = 112,000$, used in the construction of a learner. There are 11 learners, a total of $11 \times 112,000 = 1,232,000$ classifiers were evaluated and constructed in our experiments.

The performance of the algorithm can be measured under the ROC curve (AUC), Smirnov statistic (K/S), Kolmogorov- F-measure (F), geometric mean (G), accuracy (Acc), and true positive rate (TPR) were calculated. The performance measures of the last four maintains the implicit classification threshold of 0.5 (i.e., if the posterior probability of positive class membership is > 0.5 , which belongs to the positive class). AUC and K/S, tests the ability of the classifier in the separation of the positive and negative classes

Table 4. NB -Sampling, $\pi < 10\%$

Level	AUC	HSD	Level	G	HSD
ROS750	0.896	A	RUS5	81.78	A
RUS25	0.896	A	SM1000	81.37	A
BSM50	0.895	A	ROS1000	80.96	A
WE-W	0.895	A	BSM1000	76.79	B
NONE	0.895	A	CBOS	76.46	B
SM50	0.895	A	OSS	70.06	C
OSS	0.894	A	WE-W	61.21	D
CBOS	0.887	A	NONE	60.72	D
Level	AUC	HSD	Level	G	HSD
SM200	0.842	A	RUS5	70.98	A
BSM50	0.841	A	WE-W	70.17	A
WE-E	0.841	A	SM1000	70.12	A
RUS90	0.840	A	BSM1000	69.99	A
ROS750	0.840	A	ROS1000	69.47	AB
NONE	0.840	A	NONE	69.23	AB
OSS	0.831	A	OSS	67.28	B
CBOS	0.825	B	CBOS	57.70	C

Table 5. C4.5N -Sampling, $\pi < 10\%$

Level	AUC	HSD	Level	G	HSD
SM100	0.886	A	RUS5	81.51	A
BSM1000	0.884	A	SM750	66.87	B
WE-E	0.882	A	ROS500	64.98	BC
ROS50	0.881	A	CBOS	64.16	BCD
RUS25	0.881	A	BSM750	63.52	BCD
NONE	0.881	A	OSS	61.97	BCD
OSS	0.856	A	WE-W	60.34	CD
CBOS	0.846	A	NONE	59.39	D
Level	AUC	HSD	Level	G	HSD
SM300	0.853	A	RUS10	76.34	A
ROS300	0.853	A	SM1000	69.74	B
BSM1000	0.844	AB	BSM1000	67.97	BC
WE-E	0.833	BC	ROS1000	64.87	CD
RUS25	0.829	BC	WE-W	62.89	CD
OSS	0.824	C	CBOS	61.58	DE
NONE	0.820	C	OSS	60.98	DE
CBOS	0.814	C	NONE	57.66	E

Table 6. LR -Sampling, $\pi < 10\%$

Level	AUC	HSD	Level	G	HSD
ROS300	0.892	A	RUS5	81.14	A
WE-W	0.890	A	CBOS	79.08	AB
NONE	0.889	A	ROS1000	77.31	AB
RUS75	0.889	A	SM1000	75.27	BC
OSS	0.888	A	BSM1000	71.45	BC
BSM50	0.887	A	OSS	56.94	D
SM50	0.886	A	WE-W	49.24	E
CBOS	0.860	B	NONE	47.54	E
Level	AUC	HSD	Level	G	HSD
ROS500	0.847	A	ROS1000	77.09	A
WE-W	0.846	A	SM1000	76.34	A
RUS75	0.843	A	RUS10	76.03	AB
SM300	0.841	A	BSM1000	75.04	AB
BSM500	0.840	A	CBOS	72.47	B
NONE	0.839	A	WE-W	52.89	C
OSS	0.839	A	OSS	49.35	CD
CBOS	0.809	B	NONE	46.28	D

tive class membership is > 0.5 , then the example is classified as belonging to the +ve class. The two K/S and AUC will check the ability of the classifier in separating the +ve and -ve classes.

1.5 RESULTS

1.5.1. EXPERIMENTAL DATA

The results of the individual learners of the first set, we only provide a small sampling. Based on the imbalanced data, datasets are categorized into four groups those with $\pi < 5\%$, $5\% < \pi < 10\%$, $10\% < \pi < 20\%$, and finally $20\% < \pi$ (π is the percentage of examples belonging to the minority class). the performance of sampling techniques for this categorization scheme is to capture differences given different levels of imbalance, primarily on the results from $\pi < 10\%$ for the learners SVM, RF, NB, C4.5N, and LR (Tables 2 to 6). The sampling techniques of each of these tables, as measured by AUC and G, along with statistical significance and made a test. In Tables 2 to 6, the first nine rows are the results for datasets with $\pi < 5\%$, while the second nine rows are for the datasets with $5\% < \pi < 10\%$. The performance measure for the values (either AUC or G) in Tables 2 to 6 are averaged on datasets with either $\pi < 5\%$ at the top of the table or $5\% < \pi < 10\%$ at the bottom of the table. An example from Table 2, SVM with ROS1000 obtained an average AUC of 0.898 over the 20 CV runs of the eight datasets with $\pi < 5\%$, and SVM with ROS1000 obtained an average AUC of 0.861 over the 20 CV runs of the 11 datasets with $5\% < \pi < 10\%$.

ANOVA analysis on each learner and group of datasets (Berenson et al., 1983) was constructed, where the factor was the sampling technique. Tukey's Honestly Significant Difference (HSD) test (SAS Institute, 2004) is a test based on statistical analysis which measured the mean value of the performance measure for the different sampling techniques. 95% statistical confidence (all of the statistical tests in this work use 95% confidence level) two sampling techniques are not significantly different with the same block letter.

Finally note that these tables show the parameter value for each of the seven types of sampling that achieved the optimal value. For example, from Table 2, ROS at 1000% obtained the highest average AUC (across all of the datasets with $\pi < 5\%$) of 0.898, followed by RUS at 5%. Note that based on the average AUC over all datasets with $\pi < 5\%$, ROS1000, RUS5, SM1000, and BSM1000. They are not significantly different from one another (the block letter 'A' in the HSD column) used in the SVM classifier. Moreover RUS5, SM1000, BSM1000, and CBOS are not significantly different from one another, since they have the block letter 'B' in the HSD column. From these five learners we present the results and two performance measures were produced due to space limitations. AUC and G were measured in this way that is threshold dependent (G) and one that is not (AUC). Observe that in this sampling, for example AUC obtained by NB are not significantly improved Table 4, either BSM, RUS, or SM which slightly provides improvement in G.

The accuracy measurement of K/S, G, F and AUC of each learners is shown in Tables 7 to 10 present the sampling technique which results in group of imbalance. The sampling technique (with 95% statistical confidence) which was significantly maintained for no sampling and it is underlined.

Table 7. Learner- AUC-Best Sampling Technique

AUC	<5%	5%-10%	10%-20%	>20%
C4.5D	<u>RUS5</u>	<u>RUS10</u>	RUS25	RUS50
C4.5N	SM100	<u>SM300</u>	WE-W	WE-W
LR	ROS300	ROS500	ROS500	NONE
MLP	RUS10	<u>ROS300</u>	ROS300	ROS200
NB	ROS750	SM200	SM750	NONE
RBF	BSM500	<u>RUS10</u>	RUS90	WE-W
RF	<u>RUS5</u>	<u>RUS10</u>	WE-W	WE-W
RIPPER	<u>RUS5</u>	<u>RUS10</u>	<u>SM750</u>	<u>SM200</u>
SVM	<u>ROS1000</u>	<u>ROS1000</u>	<u>SM200</u>	ROS100
2NN	WE-W	SM200	WE-W	WE-W
5NN	<u>BSM300</u>	SM1000	WE-W	WE-W

Table 8. Learner- G-Sampling Technique

G	<5%	5%-10%	10%-20%	>20%
C4.5D	<u>RUS5</u>	<u>RUS10</u>	<u>RUS25</u>	<u>RUS50</u>
C4.5N	<u>RUS5</u>	<u>RUS10</u>	<u>RUS25</u>	<u>RUS50</u>
LR	<u>RUS5</u>	<u>ROS1000</u>	<u>SM500</u>	<u>ROS200</u>
MLP	<u>RUS5</u>	<u>ROS1000</u>	<u>ROS500</u>	<u>ROS200</u>
NB	<u>RUS5</u>	RUS5	BSM1000	<u>BSM200</u>
RBF	<u>RUS5</u>	<u>RUS10</u>	<u>RUS25</u>	<u>ROS200</u>
RF	<u>RUS5</u>	<u>RUS10</u>	<u>RUS25</u>	<u>SM1000</u>
RIPPER	<u>RUS5</u>	<u>RUS10</u>	<u>SM750</u>	<u>SM300</u>
SVM	<u>RUS5</u>	<u>ROS1000</u>	<u>ROS500</u>	<u>ROS200</u>
2NN	<u>RUS5</u>	<u>RUS10</u>	<u>ROS200</u>	<u>ROS200</u>
5NN	<u>RUS5</u>	<u>ROS500</u>	<u>BSM1000</u>	<u>SM200</u>

Table 9. Learner- K/S-Sampling Technique

K/S	<5%	5%-10%	10%-20%	>20%
C4.5D	<u>RUS5</u>	<u>RUS10</u>	RUS25	BSM50
C4.5N	SM500	<u>SM300</u>	BSM50	WE-W
LR	ROS500	ROS1000	ROS1000	OSS
MLP	RUS10	<u>ROS1000</u>	ROS300	ROS200
NB	WE-W	WE-E	BSM50	WE-W
RBF	RUS5	RUS10	RUS90	WE-W
RF	<u>RUS10</u>	SM1000	WE-W	WE-W
RIPPER	<u>RUS5</u>	<u>RUS10</u>	<u>SM750</u>	<u>SM300</u>
SVM	<u>ROS1000</u>	<u>ROS1000</u>	<u>ROS300</u>	ROS100
2NN	WE-W	SM1000	WE-W	WE-W
5NN	ROS750	SM1000	WE-E	WE-W

Table 10. Learner- F-Sampling Technique

F	<5%	5%-10%	10%-20%	>20%
C4.5D	<u>SM300</u>	<u>SM300</u>	<u>SM100</u>	<u>RUS50</u>
C4.5N	SM200	SM300	<u>SM100</u>	<u>WE-W</u>
LR	<u>ROS300</u>	<u>ROS500</u>	<u>SM300</u>	<u>ROS200</u>
MLP	ROS300	<u>ROS300</u>	<u>SM200</u>	<u>ROS200</u>
NB	<u>RUS25</u>	NONE	WE-W	<u>ROS200</u>
RBF	<u>RUS25</u>	<u>RUS25</u>	<u>SM200</u>	<u>ROS300</u>
RF	<u>SM1000</u>	SM750	<u>WE-E</u>	<u>WE-E</u>
RIPPER	<u>CBOS</u>	<u>SM1000</u>	<u>SM500</u>	<u>SM300</u>
SVM	<u>ROS300</u>	<u>SM500</u>	<u>SM300</u>	<u>ROS200</u>
2NN	<u>ROS200</u>	ROS750	WE-W	<u>WE-W</u>
5NN	<u>ROS200</u>	<u>ROS200</u>	<u>BSM100</u>	<u>SM200</u>

(G) and one that is not (AUC). It is shown the by applying SM,BSM and RUS sampling can be improved in NB as shown in table 4,Where as applying either RUS, SM, or BSM does significantly improve G.

Tables 7 to 10 present the sampling technique which re-sults in the best AUC, G, K/S, and F measures for each learner and group of imbalance. If applying the sampling technique resulted in performance that was significantly better (with 95% statistical confidence) than that of us-ing no sampling, then the technique is underlined.

Table 11 and figure 1, presents, over all 35 datasets, 11 learners, and six performance measures (AUC, K/S, G, F, Acc, and TPR), the number of times the rank of the sampling technique was 1, 2, . . . , 8. A rank of one means that the sampling technique, for a given dataset, learner, and performance measure, resulted in the highest value for the performance measure¹. RUS resulted in the best performance 748 times (or 32.0% = 748/2340), followed by ROS (408 times). OSS and CBOS were rarely the best technique (66 or 2.8% for OSS and 86 or 3.7% for CBOS). Further CBOS resulted in the worst perfor-mance (rank 8, last column) 965 or 42.0% of the time, followed by no sampling, which was the worst 862 or 37.5% of the time.

Tables 12,13 and figures 2, 3 display the ranking of each sampling technique separately for the four groups of imbalance ($\pi < 5\%$ at the top of Table 12 and $5\% < \pi < 10\%$ at the bottom, with $10\% < \pi < 20\%$ at the top of Table 13 and $\pi > 20\%$ at the bottom). Note that adding the individual cells of Tables 12 and 13 produces Table 11. Finally, Tables 14 to 16(figures 4,5,and 6) show the rankings of the sampling techniques only for datasets with $\pi < 5\%$ and sep-arately for each of the six performance measures, AUC, K/S, G, F, Acc, and TPR (adding the individual cells of Tables 14 to 16 produces the top half of Table 12).

1.5.2. DISCUSSION OF RESULTS

Based on the experiments conducted in this work, a number of conclusions can be drawn. The utility of sampling depends on numerous factors. First, different types of sampling work best with different learners. RUS worked very well for C4.5D (not shown) and RF, while ROS works well with LR. Second, the value of sampling is heavily dependent on the performance measure being

used. AUC and K/S, which are classification-threshold independent, generate different results than G, F, TPR, and Acc, which utilize the standard 0.5 threshold on the posterior probability. For numerous learners, such as NB, LR, 2NN, and 5NN (and to a slightly lesser extent RBF and MLP), none of the sampling techniques significantly improved the performance of the learner as measured by the AUC or K/S. However, when the performance is measured using the threshold dependent measures, significant improvements for all learners are

method	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	274	352	470	451	250	246	209	58
CBOS	86	112	115	170	406	180	276	965
NONE	236	130	147	115	165	248	407	862
OSS	66	135	167	128	234	482	809	289
ROS	408	442	365	410	325	209	145	6
RUS	748	354	367	369	270	118	67	17
SM	302	586	488	362	249	208	97	18
WE	220	195	184	410	410	610	306	82

Table 11
Datasets-. Rank of
Sampling

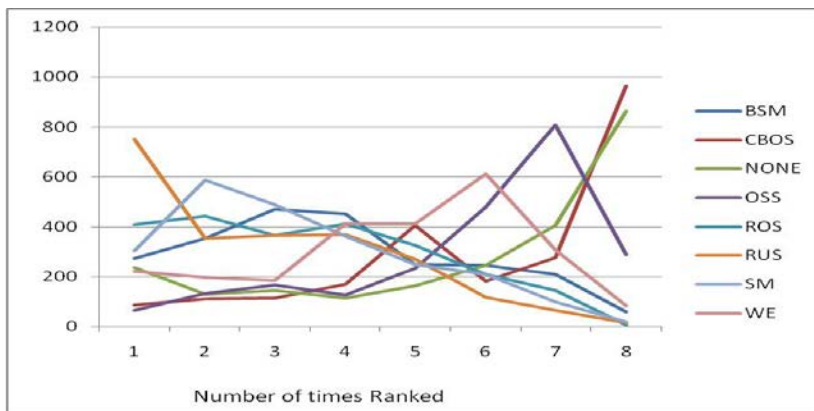


Figure 1: Shows the distribution of rank sampling

Rank based sampling is applied to various methods from the graph we identify the imbalance class

IJSER

Table 12. Datasets $\pi < 10\%$ -Rank of Sampling Techniques

Method	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	48	68	90	72	81	93	54	2
CBOS	45	59	38	30	101	26	58	171
NONE	44	40	30	25	40	82	94	173
OSS	22	42	35	29	56	84	142	118
ROS	107	91	94	120	63	31	19	3
RUS	212	89	93	61	42	12	17	2
SM	37	104	99	118	76	57	31	6
WE	18	39	44	70	73	140	117	27

Method	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	86	112	177	151	69	59	62	10
CBOS	22	23	35	73	140	74	66	293
NONE	84	27	30	37	50	82	123	286
OSS	26	37	35	36	57	128	294	103
ROS	107	133	94	143	105	72	40	3
RUS	273	99	93	103	94	46	10	4
SM	113	227	99	90	75	50	29	4
WE	39	61	44	99	131	211	104	16

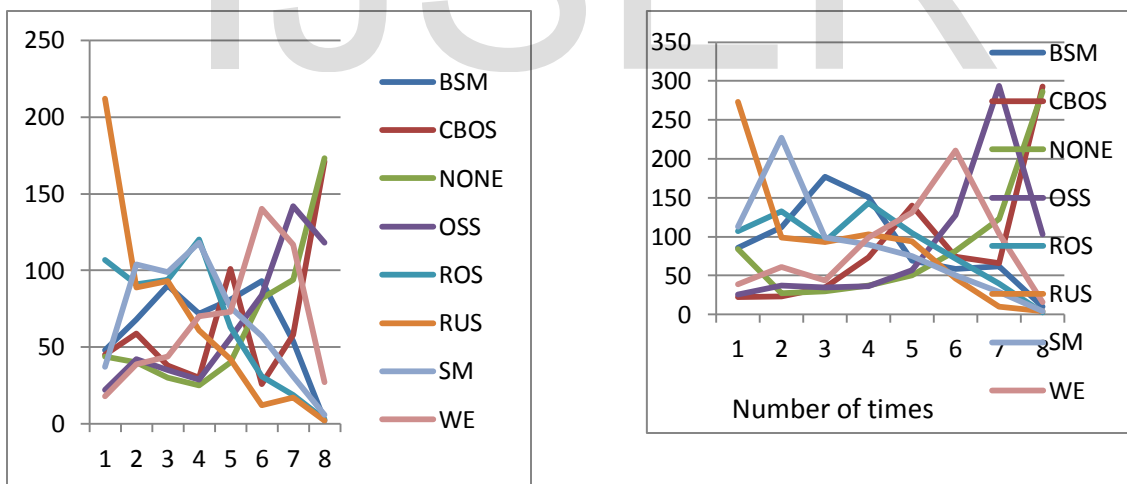


Figure 3: Shows the rank based sampling for 10%

obtained. For NB, for example, none of the sampling techniques improved the performance on datasets with $\pi < 5\%$ as measured by the AUC, however, relative to G, RUS, SM, and ROS significantly improved the performance (RUS, SM, and ROS achieved $G > 80$, while no sampling resulted in $G = 60.72$).

Table 13. Datasets $\pi > 10\%$ -Rank of Sampling Techniques,

Method	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	72	116	75	93	33	22	27	18
CBOS	14	17	13	33	56	25	66	172
NONE	33	17	32	19	35	37	77	146
OSS	6	20	26	25	40	99	136	44
ROS	70	82	55	52	55	33	29	0
RUS	104	48	61	86	56	28	12	1
SM	55	111	105	49	35	25	15	1
WE	43	44	33	37	86	105	34	14

Method	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	68	116	128	135	67	72	66	8
CBOS	5	13	29	34	109	55	86	329
NONE	75	46	48	34	40	47	113	257
OSS	12	36	61	38	81	171	237	24
ROS	124	136	93	95	102	53	57	0
RUS	159	118	116	119	78	32	28	10
SM	97	144	146	105	63	76	22	7
WE	120	51	42	97	120	154	51	25

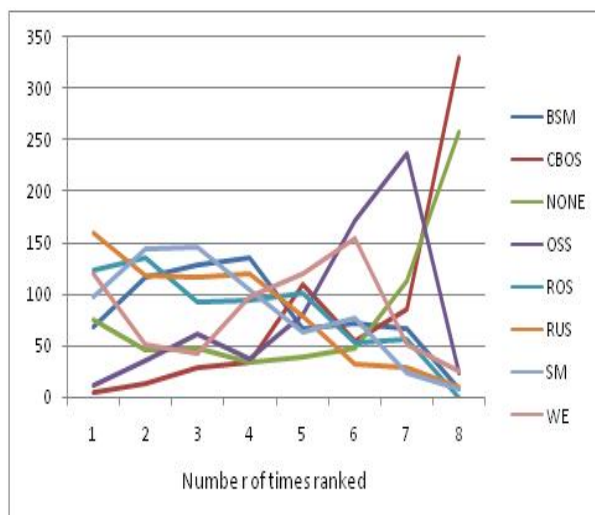


Figure 4: Shows the rank based sampling for >10%

Table 14. Datasets $\pi < 5\%$, AUC and K/S- Rank of Sampling Techniques

AUC	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	13	15	12	12	7	14	8	7
CBOS	2	8	3	2	11	6	15	41
NONE	1	8	10	8	19	20	8	14
OSS	3	8	4	8	14	14	22	22
ROS	22	1	16	19	10	6	3	0
RUS	35	12	16	10	7	3	5	0
SM	7	23	12	11	8	14	13	0
WE	5	11	13	18	13	10	4	4

K/S	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	16	18	15	5	97	13	6	6
CBOS	9	4	6	4	12	1	14	8
NONE	1	8	10	7	8	24	16	14
OSS	4	3	5	4	10	16	22	24
ROS	22	13	18	15	14	4	2	0
RUS	26	115	12	15	14	0	6	0
SM	8	21	12	18	5	13	9	2
WE	5	8	10	17	20	14	12	4

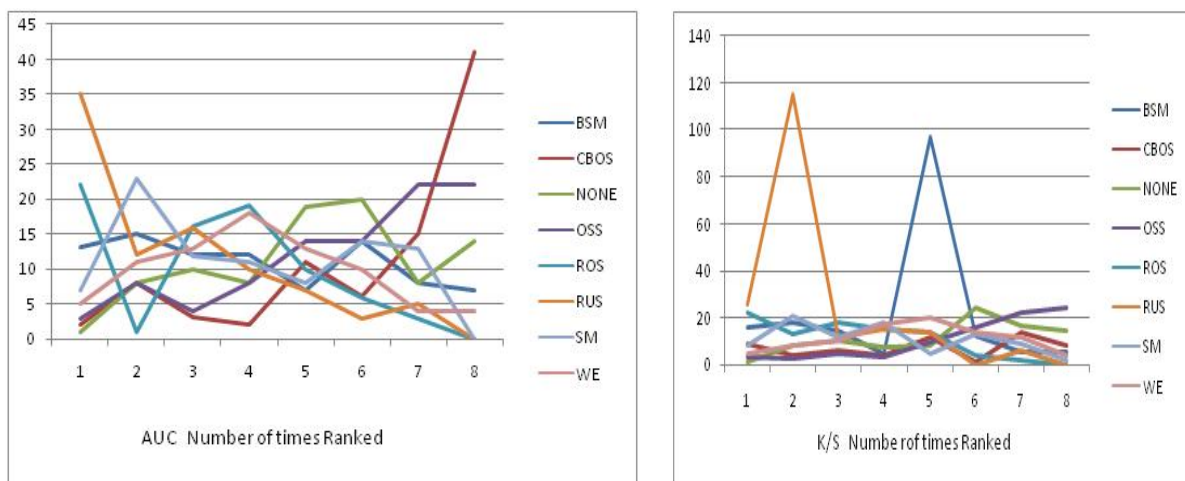


Figure 5: Shows the rank based sampling for $>5\%$ and K/S- Rank of Sampling Techniques

For AUC, with unaltered data the performance can be improved significantly in 15 of 44 scenarios (12 of the 15 occurrences with $\pi < 10\%$). For K/S, the performance of sampling is improved in 12 of the 44 scenarios. For G and F, however, in 42 and 34 of 44 scenarios, respectively, it should be outperformed without sampling. Compare to the other individual learners or datasets, RUS are best and well set in our experiments. In some cases other methods were also aimed better. In our evaluation, RUS plays a good role totally 748 of 2340 times sampling techniques. The second best was ROS, followed by BSM and SM. OSS and CBOS are very poor compared to the remaining sampling techniques, 965 of 2297 times which was shown in Tables 12 and 13, where RUS gives 39.8% and 36.4% of the time for datasets with $\pi < 5\%$ and $5\% < \pi < 10\%$. From Tables 14 to 16, RUS maintains a slight variation on ROS as the second best sampling technique compared to the AUC, K/S, and F. based on the detection of positive class (Table 16), are interested in detecting examples misleading this measure. For imbalanced data, we assume that no measure gives overall accuracy in our work. TPR, RUS is most successful when considered.

Knowledge based sampling techniques BSM, WE, SM, CBOS and OSS (mainly OSS and CBOS) Which are used to identify the inferior calculation from experiments? Which are cross checked using statistical analysis. It shows the best solutions, that are given by F and G from the validation results

Table 15. Rank of Sampling Techniques, Datasets $\pi < 5\%$,

G	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	6	8	19	14	19	17	2	3
CBOS	6	14	15	11	28	6	1	7
NONE	0	0	1	0	4	10	22	51
OSS	3	7	4	3	6	17	30	18
ROS	16	19	21	28	4	0	0	0
RUS	53	20	9	2	4	0	0	0
SM	5	19	19	27	16	0	2	0
WE	0	10	0	3	7	38	33	7

F	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	8	17	18	16	11	10	8	3
CBOS	9	5	1	2	27	10	4	5
NONE	1	5	2	5	4	17	17	60
OSS	4	8	4	4	8	16	35	10
ROS	26	20	11	21	5	4	0	0
RUS	28	11	15	17	10	3	0	0
SM	10	18	32	16	6	4	1	0
WE	3	3	5	7	18	23	32	07

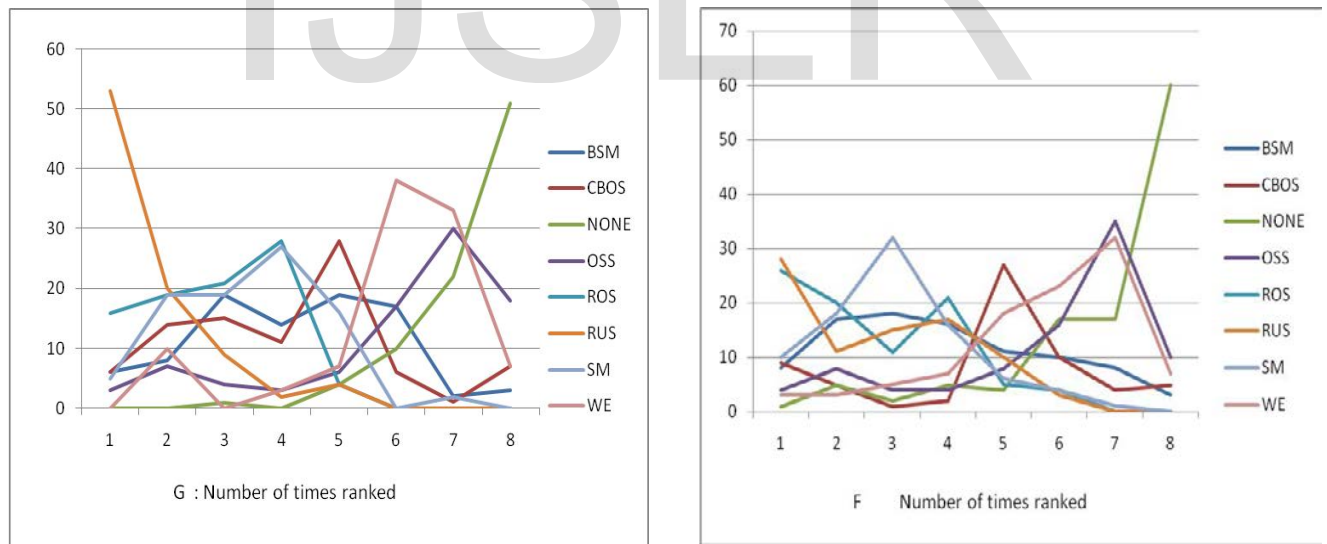


Figure 6: Shows the rank based sampling for $<5\%$ and G & C- Rank of Sampling Techniques

Table 16. Rank of Sampling Techniques, Datasets $\pi < 5\%$, Acc and T P R

ACC	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	2	5	3	12	16	19	28	3
CBOS	7	7	1	1	0	2	12	58
NONE	41	19	7	4	1	5	11	0
OSS	7	15	14	7	9	3	13	20
ROS	17	5	3	9	22	16	13	3
RUS	6	15	34	16	7	6	3	1
SM	3	4	10	15	23	25	5	3
WE	7	17	16	24	8	14	2	0
TPR	Number of times ranked							
	1	2	3	4	5	6	7	8
BSM	3	5	23	13	19	20	2	3
CBOS	12	21	12	10	23	1	4	5
NONE	0	0	0	1	4	6	17	60
OSS	1	8	4	3	9	18	35	10
ROS	4	22	25	28	8	1	0	0
RUS	64	16	7	1	0	0	0	0
SM	4	19	14	31	18	1	1	0
WE	0	0	0	1	7	41	32	07

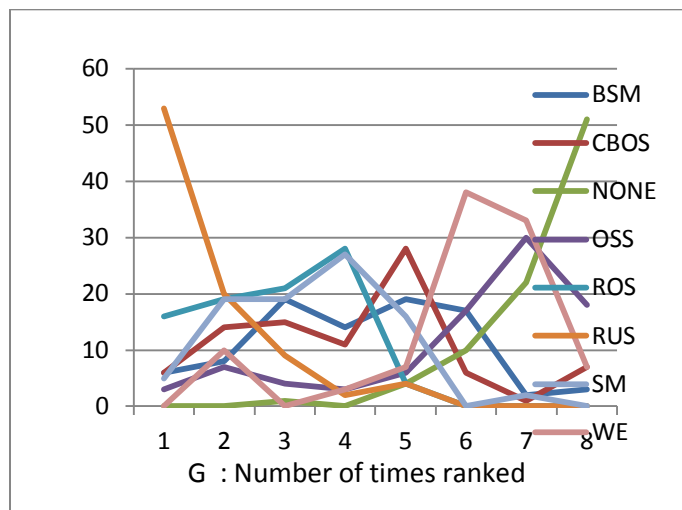
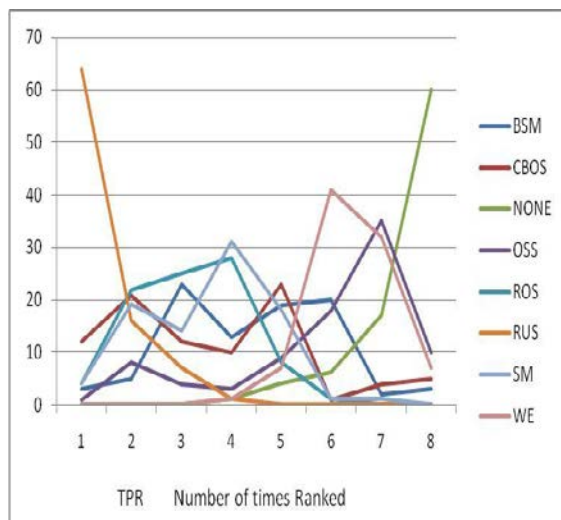


Figure 5: Shows the rank based sampling for $>5\%$ and Acc and T P R Rank of Sampling Techniques

In our Experiments, we evaluate RUS or ROS are best sampling techniques and moreover we also evaluate CBOS and OSS are very poor and gives the worst results, In very rare cases those are best sampling techniques

1.5.3. VALIDITION

There are two different types of threats the internal validity and the external validity, the internal validity which should not make any influence on the results and the external validity maintains the generalization and it tends the experimental settings which influence the outside results

In machine learning research, the experiments conducted by using WEKA. Some sampling techniques were tested thoroughly which was conducted by ANOVA analysis using SAS GLM procedure. This ANOVA analysis gives 100% reliability in results which were crossly verified by taking 35 real world datasets reduces the anomalous results. Verification can be done over one million learners. This can be mainly useful to give our conclusion reliably

Four different sampling techniques ROS, SM,ROS and BSM uses the 'free' parameter. Various possibilities and estimations can be taken to optimize the sampling percentage. By using cross validation user can estimate a value which is mainly used for the best and reliable results. Compared to other choices sampling technique is better and gives balanced results for our experiments. Dramatically these results did not change the sampling percentage. For example, among the four techniques ,RUS5 was the best technique for C4.5D for the datasets with $\pi < 5\%$ with respect to AUC (Table 7), Moreover, with CBOS and OSS, add/remove instances can be described explicitly, Therefore, a free parameter was unbiased in the comparison of sampling technique.

1.6. CONCLUSION

We give a brief , a systematical and comprehensive testing analysis of learning using imbalanced data, by using eleven algorithms of learning with 35 real word datasets from various applicational domains. The objective of the work is to guide the research to learn and practise machine learning algorithms and build the classifiers from class imbalanced datasets, and also to give a brief directions to the researchers for future research. To the best of my knowledge , no back ground work is available for analysis of imbalanced datasets with a scope, and also to compare learners using various sampling measures and check the performance of learners using different datasets. We have clearly shown that sampling is difficult to improve the performance of classifiers, mainly in geometric mean. Basic learners have responded differently in various applications. Our work is mainly based on decision tree learning, however these results show that the observations made for decision trees will not carry over to neural networks, regression, or nearest neighbor classification algorithms. Future work may consider additional learners, e.g., different variations of neural network or SVM learners. Sampling can also be compared to cost sensitive learning in future work. Alternative measures of classifier performance can also be analyzed. Future work should also consider sampling in the context of multi-class learning.

REFERENCES

- [1]. C.V. Krishna Veni, T. Sobha Rani —On the Classification of Imbalanced Datasets|| IJCST Vol . , SP 1, December 2011
- [2]. Nitesh V. Chawla, Nathalie Japkowicz, Aleksander Koleszczak —Editorial: Special Issue on Learning from Imbalanced Data Sets|| SIGKDD Explorations. Volume 6, Issue 1
- [3]. Hulse, J., Khoshgoftaar, T., Napolitano, A —Experimental perspectives on learning from im-balanced data, In: Proceedings of the 24th International Conference on Machine learning, pp. 935–942 (2007)
- [4]. Chawla, N., Bowyer, K., Hall, L., Kegelmeyer, W. —SMOTE: Synthetic minority over-sampling technique. Journal of Artificial Intelligence Research 16, 321–357 (2002)
- [5]. Andrew Estabrooks, Taeho Jo and Nathalie Japkowicz —Multiple Resampling Method for Learning from Imbalanced Data Sets. Computational Intelligence 20 (1) (2004) 18-36