

# A Comparative Study of Filter-based Feature Ranking Techniques

Huanjing Wang  
Western Kentucky University  
Bowling Green, Kentucky 42101  
huanjing.wang@wku.edu

Taghi M. Khoshgoftaar  
Florida Atlantic University  
Boca Raton, Florida 33431  
taghi@cse.fau.edu

Kehan Gao  
Eastern Connecticut State University  
Willimantic, Connecticut 06226  
gaok@easternct.edu

**Abstract**—One factor that affects the success of machine learning is the presence of irrelevant or redundant information in the training data set. Filter-based feature ranking techniques (rankers) rank the features according to their relevance to the target attribute and we choose the most relevant features to build classification models subsequently. In order to evaluate the effectiveness of different feature ranking techniques, a commonly used method is to assess the classification performance of models built with the respective selected feature subsets in terms of a given performance metric (e.g., classification accuracy or misclassification rate). Since a given performance metric usually can capture only one specific aspect of the classification performance, it may be unable to evaluate the classification performance from different perspectives. Also, there is no general consensus among researchers and practitioners regarding which performance metrics should be used for evaluating classification performance. In this study, we investigated six filter-based feature ranking techniques and built classification models using five different classifiers. The models were evaluated using eight different performance metrics. All experiments were conducted on four imbalanced data sets from a telecommunications software system. The experimental results demonstrate that the choice of a performance metric may significantly influence the classification evaluation conclusion. For example, one ranker may outperform another when using a given performance metric, but for a different performance metric the results may be reversed. In this study, we have found five distinct patterns when utilizing eight performance metrics to order six feature selection techniques.

## I. INTRODUCTION

Developing high-quality software is an important goal for any development team. Software metrics (features) that are collected during the software development process include valuable information about a software project's status, progress, quality, and evolution. Predicting the quality of software modules using software metrics in the early stages of the software development process is very critical. However, not all software metrics are relevant to the class attribute. Feature selection [1] is a process of selecting a subset of relevant features for building learning models. When irrelevant features are eliminated from the original data set, the predictive accuracy of quality models can be improved [2]. The quality models are evaluated based on performance metrics computed after the model-training process. Generally, a given performance metric can reflect a specific aspect of classification performance but cannot cover all the characteristics of it. In addition, the related literature lacks general agreement on which performance metrics should be used for evaluating classification performance [3], [4], [5].

In this empirical study, we investigated six different filter-based feature ranking techniques (rankers), chi-square (CS), information gain (IG), gain ratio (GR), symmetrical uncertainty (SU), and two forms of ReliefF (RF and RFW). In order to evaluate the effectiveness of these methods, we built classification models using five different classifiers on the smaller subsets of selected attributes. The five classifiers used in the study include naïve Bayes (NB), multilayer per-

ceptron (MLP),  $k$ -nearest neighbors (KNN), support vector machine (SVM), and logistic regression (LR). Each classification model is assessed with eight different performance metrics: the area under the Receiver Operating Characteristic (ROC) curve (AUC), the area under the Precision-Recall curve (PRC), Default F-Measure (DFM), Best F-Measure (BFM), Default Geometric Mean (DGM), Best Geometric Mean (BGM), Default Arithmetic Mean (DAM), and Best Arithmetic Mean (BAM).

The empirical validation of the different models was implemented through a case study of four imbalanced data sets from a telecommunications software system. Each data set holds the same number of attributes but has a different number of observations. The results demonstrate that the selection of a performance metric may directly impact the evaluation outcome. For instance, one ranker may perform better than another ranker in terms of a given performance metric, but this may not be true when using a different performance metric. In this study, we have discovered five distinct patterns when we used eight performance metrics to order six feature selection techniques.

The main contribution of this work is to provide an assessment and comparison of six filter-based feature ranking techniques using eight performance metrics and over five different classifiers. To our knowledge, no one has done such an extensive study yet.

The rest of the paper is organized as follows. Section II provides more detailed information about the techniques used in the study. The software measurement data sets used in the experiment are described in Section III. Section IV presents the experimental results and analysis. Finally, the conclusion is summarized in Section V.

## II. METHODOLOGY

### A. Filter-based Feature Ranking Techniques

Filter-based feature ranking techniques rank features independently without involving any learning algorithm. Feature ranking consists of scoring each feature according to a particular method, then selecting features based on their scores. This work employs some commonly used filter-based feature ranking techniques including chi-square, information gain, gain ratio, symmetrical uncertainty, and ReliefF. The chi-square (CS) [6] test is used to examine if there is 'no association' between two attributes, i.e., whether the two variables are independent. Information gain, gain ratio, and symmetrical uncertainty are measures based on the concept of entropy, which is based on information theory. Information gain (IG) [7] is the information provided about the target class attribute  $Y$ , given the value of independent attribute  $X$ . Information gain measures the decrease of the weighted average impurity of the partitions, compared with the impurity of the complete set of data. A drawback of IG is that it tends to prefer attributes with a larger number of possible values. One strategy to counter this problem is to use the gain ratio (GR), which

penalizes multi-valued attributes. Symmetrical uncertainty (SU) [8] is another way to overcome the problem of IG’s bias toward attributes with more values, doing so by dividing IG by the sum of the entropies of X and Y. Relief is an instance-based feature ranking technique [9]. ReliefF is an extension of the Relief algorithm that can handle noise and multi-class data sets. When the ‘weightByDistance’ (weight nearest neighbors by their distance) parameter is set as default (false), the algorithm is referred to as RF; when the parameter is set to true, the algorithm is referred to as RFW.

### B. Classifiers

Software quality models are built with five well-known classification algorithms [10] including naïve Bayes (NB), multilayer perceptron (MLP),  $k$ -nearest neighbors (KNN), support vector machine (SVM) and logistic regression (LR). These were selected because of their common use in software engineering and other data mining applications. Unless stated otherwise, we use default parameter settings for the different learners as specified in the WEKA [10] data mining tool. Parameter settings are changed only when a significant improvement in performance is obtained. For the KNN classifier, 5 neighbors are used in the study.

### C. Performance Metrics

In a two-group classification problem, such as fault-prone (positive) and not fault-prone (negative), there can be four possible prediction outcomes: true positive (TP) (i.e., correctly classified positive instances), false positive (FP) (i.e., negative instance classified as positive), true negative (TN) (i.e., correctly classified as negative instance), and false negative (FN) (i.e., positive instance classified as negative). The numbers of cases from the four sets (outcomes) form the basis for several other performance measures that are well known and commonly used for classifier evaluation.

- *Area Under ROC (Receiver Operating Characteristic) Curve (AUC)*: has been widely used to measure classification model performance [11]. AUC is a single-value measurement that ranges from 0 to 1. The ROC curve is used to characterize the trade-off between true positive rate ( $\frac{|TP|}{|TP|+|FN|}$ ) and false positive rate ( $\frac{|FP|}{|FP|+|TN|}$ ). A perfect classifier provides an AUC that equals 1.
- *Area Under the Precision-Recall Curve (PRC)*: is a single-value measure that originated from the area of information retrieval. The area under the PRC ranges from 0 to 1. The PRC diagram depicts the trade off between recall ( $\frac{|TP|}{|TP|+|FN|}$ ) and precision ( $\frac{|TP|}{|TP|+|FP|}$ ). A classifier that is near optimal in AUC space may not be optimal in precision/recall space.
- *Default F-measure (DFM)*: The F-measure is a single value metric that originated from the field of information retrieval [12]. It is calculated as  $\frac{2|TP|}{2|TP|+|FP|+|FN|}$ . The Default F-measure (DFM) corresponds to a decision threshold value of 0.5.
- *Best F-Measure (BFM)*: is the largest value of F-measure when varying the decision threshold value between 0 and 1. A perfect classifier yields an F-measure of 1, i.e., no misclassification.
- *Default Geometric Mean (DGM)*: The Geometric Mean (GM) is a single-value performance measure that ranges from 0 to 1, and a perfect classifier provides a value of 1. GM is defined as the square root of the product of true positive rate and true negative rate, where the true negative rate is defined as  $\frac{|TN|}{|FP|+|TN|}$ . The decision threshold  $t = 0.5$  is used for the Default Geometric Mean (DGM).

TABLE I  
SOFTWARE DATA SETS CHARACTERISTICS

	Data	#Metrics	#Modules	%fp	%nfp
LLTS	SP1	42	3649	6.28%	93.72%
	SP2	42	3981	4.75%	95.25%
	SP3	42	3541	1.33%	98.67%
	SP4	42	3978	2.31%	97.69%

TABLE II  
PERFORMANCE METRICS USING NB

Data	Ranker	AUC	PRC	DFM	BFM	DGM	BGM	DAM	BAM
SP1	CS	0.7846	0.2331	0.2895	0.3045	0.5616	0.7227	0.6356	0.7241
	IG	0.7346	0.216	0.2777	0.2942	0.5113	0.6966	0.6138	0.7018
	GR	0.7831	0.2271	0.294	0.3109	<b>0.5706</b>	0.7204	0.6404	0.7220
	RF	0.7879	0.213	0.2706	0.2953	0.5394	0.7301	0.6226	0.7309
	RFW	<b>0.7882</b>	0.2145	0.2682	0.2888	0.5397	<b>0.7320</b>	0.6222	<b>0.7326</b>
	SU	0.7865	<b>0.2420</b>	<b>0.3046</b>	<b>0.3140</b>	0.5676	0.7214	<b>0.6411</b>	0.7231
SP2	CS	<b>0.8108</b>	0.1975	<b>0.2797</b>	0.2915	<b>0.5891</b>	<b>0.7526</b>	<b>0.6535</b>	<b>0.7532</b>
	IG	0.7613	<b>0.1988</b>	0.2793	<b>0.2967</b>	0.5314	0.7217	0.627	0.7241
	GR	0.8081	0.1886	0.2629	0.272	0.5617	0.7524	0.6376	0.7528
	RF	0.8053	0.1941	0.2409	0.2649	0.5353	0.7314	0.622	0.7337
	RFW	0.8081	0.1974	0.242	0.2677	0.5367	0.7335	0.6228	0.7363
	SU	0.7729	0.1806	0.2682	0.2831	0.5511	0.7281	0.6341	0.7295
SP3	CS	0.8184	0.072	0.1319	0.1561	0.5663	0.7689	0.6435	0.7705
	IG	0.7808	0.0603	0.1203	0.1398	0.53	0.7437	0.6234	0.7457
	GR	0.8118	0.0721	<b>0.1384</b>	0.1563	<b>0.5884</b>	0.7663	<b>0.6566</b>	0.7678
	RF	<b>0.8305</b>	<b>0.0767</b>	0.1285	<b>0.1608</b>	0.5435	<b>0.7952</b>	0.6316	<b>0.7957</b>
	RFW	0.819	0.0744	0.1303	0.1596	0.5492	0.7662	0.6346	0.7688
	SU	0.7882	0.0645	0.1238	0.1467	0.5446	0.7476	0.6311	0.7489
SP4	CS	0.7696	0.1229	0.2094	0.2358	<b>0.6211</b>	0.7286	<b>0.6757</b>	0.7328
	IG	0.7519	0.1121	<b>0.2189</b>	0.2307	0.5798	0.722	0.654	0.7298
	GR	<b>0.7794</b>	0.1103	0.1943	0.2098	0.5984	<b>0.7292</b>	0.6605	<b>0.7332</b>
	RF	0.7731	0.124	0.2146	0.245	0.5967	0.7267	0.6624	0.7295
	RFW	0.7735	<b>0.1273</b>	0.2172	<b>0.2533</b>	0.6002	0.726	0.6646	0.7286
	SU	0.7592	0.1105	0.2124	0.2257	0.5891	0.7262	0.658	0.7317

- *Best Geometric Mean (BGM)*: is the maximum Geometric Mean value that is obtained when varying the decision threshold between 0 and 1.
- *Default Arithmetic Mean (DAM)*: The arithmetic mean is just like the geometric mean but uses the arithmetic mean of the true positive rate and true negative rate instead of the geometric mean. It is also a single-value performance measure that ranges from 0 to 1. The decision threshold  $t = 0.5$  is used for the Default Arithmetic Mean (DAM).
- *Best Arithmetic Mean (BAM)*: is just like the BGM, but using the maximum arithmetic mean that is obtained when varying the decision threshold between 0 and 1.

### III. DATA SET CHARACTERISTICS

Experiments conducted in this study used software metrics and defect data collected from a real-world software project, a very large telecommunications software system (denoted as LLTS) [13]. LLTS contains data from four consecutive releases, which are labeled as SP1, SP2, SP3, and SP4. The software measurement data sets consist of 42 software metrics, including 24 product metrics, 14 process metrics, and four execution metrics [13]. The dependent variable is the class of the program module, fault-prone ( $fp$ ), or not fault-prone ( $nfp$ ). A program module with one or more faults is considered  $fp$ , and  $nfp$  otherwise. Table I lists the characteristics of the four release data sets utilized in this work. An important characteristic of these data sets is that they all suffer from class imbalance, where the proportion of  $fp$  modules is much lower than that of  $nfp$  modules.

### IV. EXPERIMENTS

#### A. Experimental Design

We first used six filter-based rankers to select the subsets of attributes. We ranked the features and selected the top  $\lceil \log_2 n \rceil$  features according to their respective scores, where  $n$  is the number of independent features for a given data set. The reasons why we select the top  $\lceil \log_2 n \rceil$  features include (1) related literature does not provide guidance on the appropriate number of features to select; and



TABLE VII  
ANALYSIS OF VARIANCE

(a) AUC					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.2329	5	0.0466	53.78	0
B	3.8618	4	0.9655	1114.88	0
A × B	0.0684	20	0.0034	3.95	0
Error	1.0132	1170	0.0009		
Total	5.1763	1199			

(b) PRC					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.0223	5	0.0045	1.08	0.370
B	1.2106	4	0.3027	73.28	0
A × B	0.0583	20	0.0029	0.71	0.823
Error	4.8325	1170	0.0041		
Total	6.1237	1199			

(c) DFM					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.0602	5	0.0120	4.39	0.001
B	5.6876	4	1.4219	519.19	0
A × B	0.0085	20	0.0004	0.16	1.000
Error	3.2043	1170	0.0027		
Total	8.9606	1199			

(d) BFM					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.0240	5	0.0048	1.41	0.216
B	1.7631	4	0.4408	129.66	0
A × B	0.0975	20	0.0049	1.43	0.097
Error	3.9773	1170	0.0034		
Total	5.8619	1199			

(e) DGM					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.1137	5	0.0227	2.55	0.026
B	35.8865	4	8.9716	1006.31	0
A × B	0.1282	20	0.0064	0.72	0.809
Error	10.4310	1170	0.0089		
Total	46.5594	1199			

(f) BGM					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.0954	5	0.0191	27.51	0
B	2.3820	4	0.5955	858.74	0
A × B	0.0699	20	0.0035	5.04	0
Error	0.8114	1170	0.0007		
Total	3.3587	1199			

(g) DAM					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.0084	5	0.0017	8.01	0
B	3.0634	4	0.7659	3645.18	0
A × B	0.0118	20	0.0006	2.81	0
Error	0.2458	1170	0.0002		
Total	3.3294	1199			

(h) BAM					
Source	Sum Sq.	d.f.	Mean Sq.	F	<i>p</i> -value
A	0.0748	5	0.0150	26.00	0
B	1.9306	4	0.4827	838.88	0
A × B	0.0824	20	0.0041	7.16	0
Error	0.6732	1170	0.0006		
Total	2.7609	1199			

filter-based rankers were considered, and Factor B, in which five classifiers were included. In addition, the interaction A×B was also included. In this ANOVA test, the results from all four release data sets were taken into account together. A significance level of  $\alpha = 5\%$  was used for all statistical tests.

The ANOVA results are presented in Table VII. From the table, we can see that for the performance metrics AUC, BGM, DAM and BAM, the *p*-values for the main factors A and B, and the interaction term A×B were zeros, indicating the performance values are not same for all groups in each of the main factors and also influenced by the interaction term A×B, i.e., Factor A is different at every level of Factor B, and vice versa. For the performance metrics PRC and BFM, there was no significant distinction between any pair of the group means for Factor A and interaction A×B since the *p*-values were greater than 0.05, while an obvious difference existed in at least one pair of group means for Factor B, because the *p*-value was zero. For the performance metrics DFM and DGM, an obvious difference existed in at least one pair of group means for Factor A and also for Factor B. However, their interaction did not contribute too much for

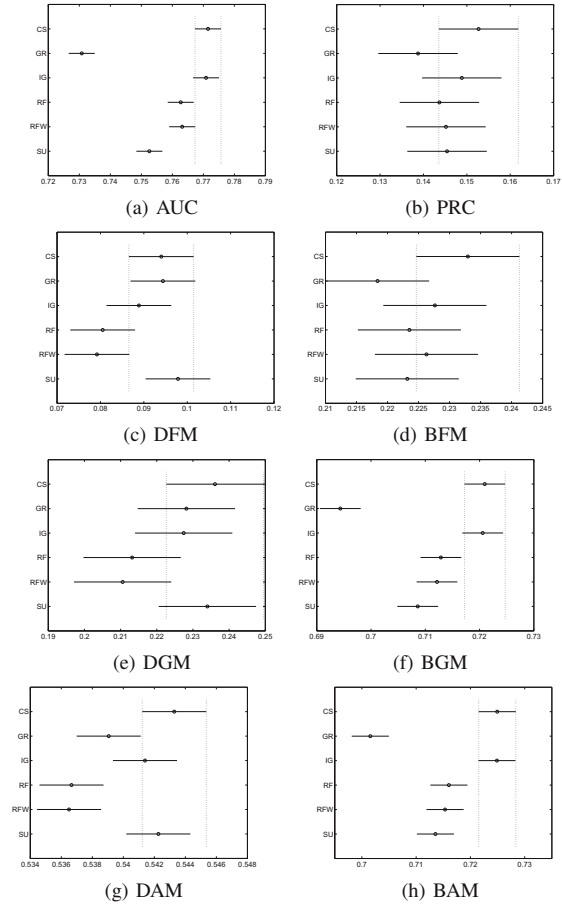


Fig. 2. Multiple comparisons for Factor A

the classification performance.

Additional multiple comparisons for the main factors and interaction term were performed to investigate the differences among the respective groups (levels). Both ANOVA and multiple comparison tests were implemented in MATLAB. The multiple comparisons are presented in Fig. 2 through 4. The performance of filter-based rankers was ranked from best to worst for each performance metric as shown in Table VIII. Each ranker is labeled with a superscript. The rankers labeled with the same superscripts implies that they were from same performance group, in which no statistically significant difference was found between rankers. The table shows five distinct groups of results when we order six commonly used rankers based on eight performance metrics (over all the classifiers built): (1) PRC, DGM, and BFM (when using these three metrics to evaluate the rankers, the orders of the six rankers are the same or similar.); (2) BGM and BAM (identical ordering of six feature-based rankers); (3) AUC; (4) DFM; and (5) DAM. The performance of learners was also ranked from best to worst for each performance metric as shown in Table IX. We can observe that three distinct patterns emerge when we are ordering learners based on eight performance metrics: (1) AUC, BGM, and BAM; (2) PRC and BFM; and (3) DFM, DGM, and DAM. All the ranks of interaction of rankers and learners are also summarized but not presented here due to space limitations.

Some findings can be summarized from these tables and figures.

- For all performance metrics, there are no significant differences

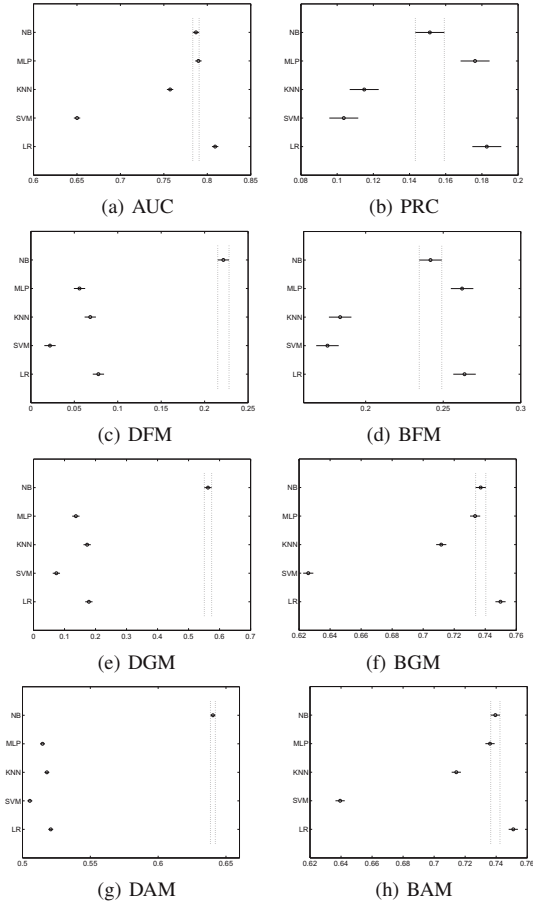


Fig. 3. Multiple comparisons for Factor B

TABLE VIII  
RANK OF RANKERS (FACTOR A)

Ranker	Ranks (best $\rightarrow$ worst)					
AUC	CS <sup>a</sup>	IG <sup>a,b</sup>	RFW <sup>a,b</sup>	RF <sup>b</sup>	SU <sup>c</sup>	GR <sup>d</sup>
PRC	CS <sup>a</sup>	IG <sup>a</sup>	SU <sup>a</sup>	RFW <sup>a</sup>	RF <sup>a</sup>	GR <sup>a</sup>
DFM	SU <sup>a</sup>	GR <sup>a,b</sup>	CS <sup>a,b,c</sup>	IG <sup>a,b,c</sup>	RF <sup>b,c</sup>	RFW <sup>c</sup>
BFM	CS <sup>a</sup>	IG <sup>a</sup>	RFW <sup>a</sup>	RF <sup>a</sup>	SU <sup>a</sup>	GR <sup>a</sup>
DGM	CS <sup>a</sup>	SU <sup>a</sup>	GR <sup>a</sup>	IG <sup>a</sup>	RF <sup>a</sup>	RFW <sup>a</sup>
BGM	CS <sup>a</sup>	IG <sup>a</sup>	RF <sup>b</sup>	RFW <sup>b</sup>	SU <sup>b</sup>	GR <sup>c</sup>
DAM	CS <sup>a</sup>	SU <sup>a,b</sup>	IG <sup>a,b</sup>	GR <sup>b,c</sup>	RF <sup>c</sup>	RFW <sup>c</sup>
BAM	CS <sup>a</sup>	IG <sup>a</sup>	RF <sup>b</sup>	RFW <sup>b</sup>	SU <sup>b</sup>	GR <sup>c</sup>

TABLE IX  
RANK OF CLASSIFIERS (FACTOR B)

Classifier	Ranks (best $\rightarrow$ worst)					
AUC	LR <sup>a</sup>	MLP <sup>b</sup>	NB <sup>b</sup>	KNN <sup>c</sup>	SVM <sup>d</sup>	
PRC	LR <sup>a</sup>	MLP <sup>a</sup>	NB <sup>b</sup>	KNN <sup>c</sup>	SVM <sup>c</sup>	
DFM	NB <sup>a</sup>	LR <sup>b</sup>	KNN <sup>b,c</sup>	MLP <sup>c</sup>	SVM <sup>d</sup>	
BFM	LR <sup>a</sup>	MLP <sup>a</sup>	NB <sup>b</sup>	KNN <sup>c</sup>	SVM <sup>c</sup>	
DGM	NB <sup>a</sup>	LR <sup>b</sup>	KNN <sup>b</sup>	MLP <sup>c</sup>	SVM <sup>d</sup>	
BGM	LR <sup>a</sup>	NB <sup>b</sup>	MLP <sup>b</sup>	KNN <sup>c</sup>	SVM <sup>d</sup>	
DAM	NB <sup>a</sup>	LR <sup>b</sup>	KNN <sup>b,c</sup>	MLP <sup>c</sup>	SVM <sup>d</sup>	
BAM	LR <sup>a</sup>	NB <sup>b</sup>	MLP <sup>b</sup>	KNN <sup>c</sup>	SVM <sup>d</sup>	

between CS and IG, the performance differences between RF and RFW are minimal.

- There are no significant differences when ordering all rankers in terms of PRC, DGM, and BFM performance metrics.
- One method being ranked at top by a given performance metric does not mean that it is also ranked at top by another performance metric, and the same for being ranked worst. For example, GR performed worse than other filter-based rankers when using AUC to evaluate classification performance (see Fig. 2(a)), while this is not true when using a different performance metric, for instance, DFM (see Fig. 2(c)).
- CS has the best performance according to all performance metrics except DFM, while SU has the best performance for DFM.
- The performance of various ranking techniques and learners shows two different patterns. One pattern is found when AUC, PRC, BFM, BGM and BAM are utilized for assessment. For Factor A (see Fig. 2), CS performed best, followed by IG; GR performed worst among the six filter-based feature ranking techniques; and RF, RFW and SU sat in between. For Factor B (see Fig. 3), LR performed best, followed by MLP and NB, then KNN, and finally SVM. The other pattern appears when DFM, DGM and DAM are used for evaluation. The pattern is that, for Factor A, RF and RFW performed worse than the other four ranking techniques; for factor B, NB significantly outperformed all other learners, followed by LR, KNN, and MLP, and finally SVM. These two patterns are also extended to interaction A×B. The two distinct patterns can be easily observed from Fig. 4.
- The performance distributions of the 30 group means are very similar when evaluated using DFM, DGM and DAM (see Fig. 4(c), 4(e) and 4(g)). The NB group performed much better than the other groups, while the performances of the remaining four groups are relatively close to each other. But still we can see that the KNN and LR groups performed better than MLP and SVM groups. Of the two inferior performance groups, SVM performed even worse. Meanwhile, the performance distributions of the 30 group means show a similar pattern when evaluated in terms of AUC, PRC, BFM, BGM and BAM (see Fig. 4(a), 4(b), 4(d), 4(f) and 4(h)). Overall, the NB, MLP and LR groups present relatively similar performances, but still we can see that the LR group performed best. These three groups outperformed the KNN and SVM groups. In fact, the SVM performed once again worst among the five learner groups. Also, one point that needs to be noted is that if we have to compare the impacts of learners and filter ranking techniques on the classification performance, we can clearly see that learners had more influence on the classification performance in this study.

## V. CONCLUSION

In this paper, we present six filter-based feature ranking techniques and evaluate their effectiveness by building five different types of classification models. Each model is assessed in terms of eight performance metrics. The experiments were conducted on four consecutive releases of a very large telecommunications system. The experimental results demonstrate that the selection of a performance metric is critical for assessing classification performance. Using different performance metrics may generate different evaluation results. We summarized five distinct patterns of the six feature ranking techniques when using the eight performance metrics. Every metric concurred on the identification of the worst learner, SVM. These results accentuate

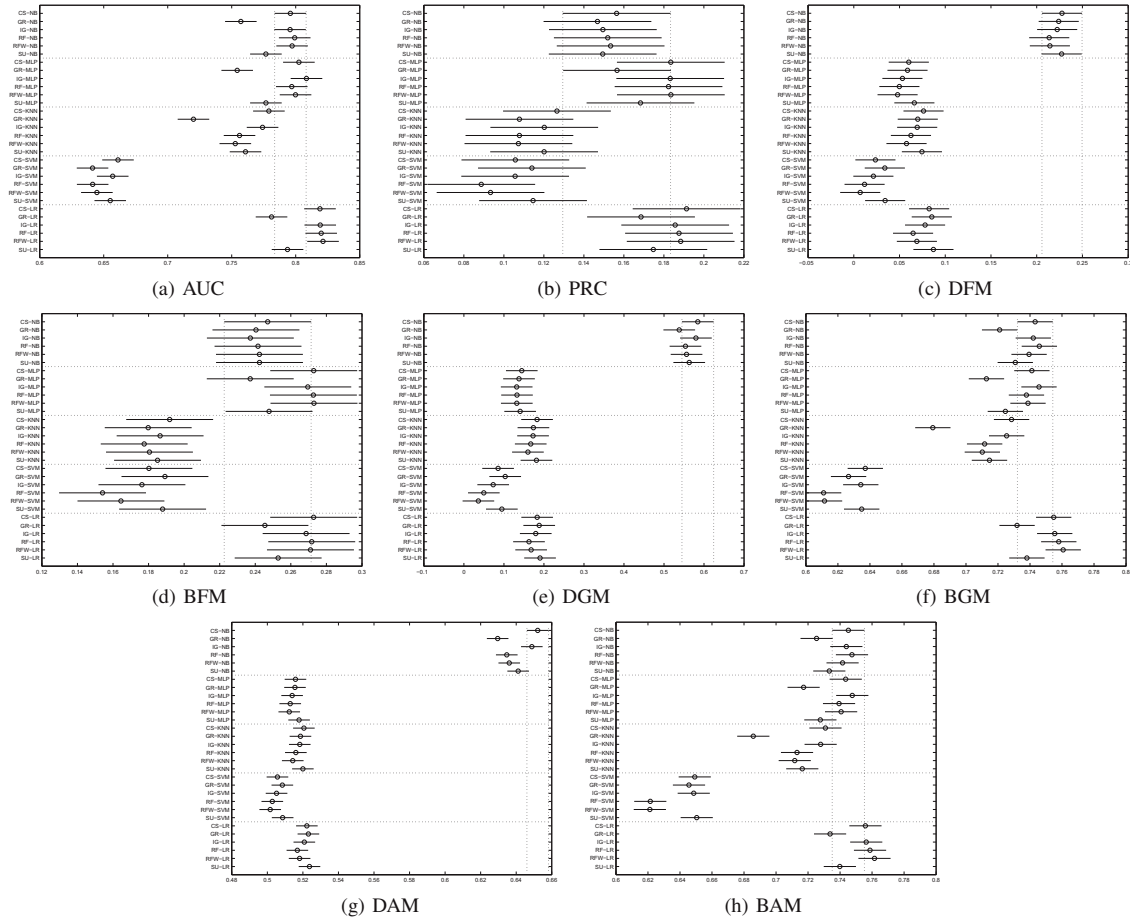


Fig. 4. Multiple comparisons for Factor A x B

the importance of metric selection for learning from class imbalanced data.

More investigations of characteristics of performance metrics and their impact on classification performance using a variety of domain data will be studied in our future work.

#### REFERENCES

- [1] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *Journal of Machine Learning Research*, vol. 3, pp. 1157–1182, March 2003.
- [2] H. Wang, T. M. Khoshgoftaar, K. Gao, and N. Seliya, "High-dimensional software engineering data and feature selection," in *Proceedings of 21st IEEE International Conference on Tools with Artificial Intelligence*, Newark, NJ, USA, Nov. 2-5 2009, pp. 83–90.
- [3] M. Sokolova, N. Japkowicz, and S. Szpakowicz, "Beyond accuracy, f-score and roc: A family of discriminant measures for performance evaluation," *AI 2006: Advances in Artificial Intelligence*, pp. 1015–1021, 2006.
- [4] N. Seliya, T. M. Khoshgoftaar, and J. Van Hulse, "A study on the relationships of classifier performance metrics," in *Proceedings of the 21st IEEE International Conference on Tools with Artificial Intelligence*, 2009, pp. 59–66.
- [5] A. Folleco, T. M. Khoshgoftaar, and A. Napolitano, "Comparison of four performance metrics for evaluating sampling techniques for low quality class-imbalanced data," in *Proceedings of the 2008 Seventh International Conference on Machine Learning and Applications*, Washington, DC, USA, 2008, pp. 153–158.
- [6] R. L. Plackett, "Karl pearson and the chi-squared test," *International Statistical Review*, vol. 51, no. 1, pp. 59–72, 1983.
- [7] J. R. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [8] M. A. Hall and G. Holmes, "Benchmarking attribute selection techniques for discrete class data mining," *IEEE Transactions on Knowledge and Data Engineering*, vol. 15, no. 6, pp. 1437 – 1447, Nov/Dec 2003.
- [9] K. Kira and L. A. Rendell, "A practical approach to feature selection," in *Proceedings of 9th International Workshop on Machine Learning*, 1992, pp. 249–256.
- [10] I. H. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques*, 2nd ed. Morgan Kaufmann, 2005.
- [11] T. Fawcett, "An introduction to ROC analysis," *Pattern Recognition Letters*, vol. 27, no. 8, pp. 861–874, June 2006.
- [12] J. Davis and M. Goadrich, "The relationship between Precision-Recall and ROC curves," in *Proceedings of the 23rd International Conference on Machine Learning*, Pittsburgh, Pennsylvania, 2006, pp. 233–240.
- [13] K. Gao, T. M. Khoshgoftaar, and H. Wang, "An empirical investigation of filter attribute selection techniques for software quality classification," in *Proceedings of the 10th IEEE International Conference on Information Reuse and Integration*, Las Vegas, Nevada, August 10-12 2009, pp. 272–277.
- [14] T. M. Khoshgoftaar, M. Golawala, and J. Van Hulse, "An empirical study of learning from imbalanced data using random forest," in *Proceedings of the 19th IEEE International Conference on Tools with Artificial Intelligence*, vol. 2, Washington, DC, USA, 2007, pp. 310–317.
- [15] M. L. Berenson, M. Goldstein, and D. Levine, *Intermediate Statistical Methods and Applications: A Computer Package Approach*, 2nd ed. Englewood Cliffs, NJ: Prentice-Hall, 1983.