

Investigating the Effect of Sampling Methods for Imbalanced Data Distributions

Show-Jane Yen, Yue-Shi Lee, Cheng-Han Lin and Jia-Ching Ying

Abstract—Classification is an important and well-known technique in the field of machine learning, and the training data will significantly influence the classification accuracy. However, the training data in real-world applications often are imbalanced class distribution. It is important to select the suitable training data for classification in the imbalanced class distribution problem. In this paper, we propose a cluster-based sampling approach for selecting the representative data as training data to improve the classification accuracy and investigate the effect of under-sampling methods in the imbalanced class distribution problem. In the experiments, we evaluate the performances for our cluster-based sampling approach and the other sampling methods in the previous studies.

I. INTRODUCTION

Classification Analysis [4, 6, 10] is a well-studied technique in data mining and machine learning domains. Due to the forecasting characteristic of classification, it has been used in a lot of real applications, such as flow-away customers and credit card fraud detections in finance corporations. Classification analysis can produce a class predicting system (or called a classifier) by analyzing the properties of a dataset with classes. The classifier can make class forecasts on new samples with unknown class labels. For example, a medical officer can use medical predicting system to predict if a patient have drug allergy or not. A dataset with given class can be used to be a training dataset, and a classifier must be trained by a training dataset to have the capability for class prediction.

The classification techniques usually assume that the training samples are uniformly-distributed between different classes. A classifier performs well when the classification technique is applied to a dataset evenly distributed among different classes. However, many datasets in real applications involve imbalanced class distribution problem [1, 2, 11, 13, 14, 19]. The imbalanced class distribution problem occurs while there are much more samples in one class than the other class in a training dataset. In an imbalanced dataset, the *majority class* has a large percent of all the samples, while the samples in *minority class* just occupy a small part of all the samples. In this case, a classifier usually tends to predict that samples have the majority class and completely ignore the minority class.

Many applications such as fraud detection, intrusion prevention, risk management, medical research often have the imbalanced class distribution problem. For example, a bank would like to construct a classifier to predict that whether the customers will have fiduciary loans in the future or not. The number of customers who have had fiduciary loans is only two percent of all customers. If a fiduciary loan classifier predicts that all the customers never have fiduciary loans, it will have a

quite high accuracy as 98 percent. However, the classifier can not find the target people who will have fiduciary loans within all customers. Therefore, if a classifier can make correct prediction on the minority class efficiently, it will be useful to help corporations make a proper policy and save a lot of cost. In this paper, we study the effects of under-sampling [19] on the neural network technique and propose some new under-sampling methods based on clustering, such that the influence of imbalanced class distribution can be decreased and the accuracy of predicting the minority class can be increased.

II. RELATED WORK

Since many real applications have the imbalanced class distribution problem, researchers have proposed several methods to solve this problem. These methods try to solve the class distribution problem both at the algorithmic level and data level. At the algorithmic level, developed methods include cost-sensitive learning [7, 8, 18] and recognition-based learning [3, 15].

Cost-sensitive learning approach assumes the misclassification costs are known in a classification problem. A cost-sensitive classifier tries to learn more characteristics of samples with the minority class by setting a high cost to the misclassification of a minority class sample. However, misclassification costs are often unknown and a cost-sensitive classifier may result in overfitting training. To ensure learning the characteristics of whole samples with the minority class, the recognition-based learning approach attempts to overfit by one-class (minority class) learning. One-class learning is more suitable than two-class approaches under certain conditions such like very imbalanced data and high dimensional noisy feature space [8].

At the data level, methods include multi-classifier committee [9, 16], and re-sampling [2, 3, 5, 7, 12, 19] approaches. Multi-classifier committee approach [9, 16] makes use of all information on a training dataset. Assume in a training dataset, MA is the sample set with majority class, and MI is the other set with minority class. Multi-classifier committee approach divides the samples with majority class (i.e. MA) randomly into several subsets, and then takes every subset and all the samples with minority class (i.e. MI) as training dataset, respectively. The number of the subsets depends on the ratio of MA's size to MI's size. For example, suppose in a dataset, the size of MA is 48 (samples) and the size of MI is 2 (samples). If we think the best ratio of MA's size to MI's size is 1:1 in a training dataset, then the number of training subsets will be $48/2=24$. Each of these 24 subsets

contains MI and a subset of MA that both sizes are 2, and the ratio of them is exactly 1:1.

After training these training datasets separately, several classifiers are available as committees. Multi-classifier committee approach uses all the classifiers to predict a sample and decides the final class to it by the prediction results of the classifiers. Voting is one simple method for making a final class decision to a sample, in which a minimum threshold is set up. If the number of classifiers that predict the same class “C” for a sample exceeds the minimum threshold, then the final class prediction of this sample will be “C”. Though multi-classifier committee approach does not abandon any sample from MA, it may be inefficient in the training time for all the committees and can not ensure the quality of every committee. Further selection of the committees will make the predictions more correct and more efficient.

As for re-sampling approach, it can be distinguished into *over-sampling approach* [3, 5, 12] and *under-sampling approach* [2, 19]. The over-sampling approach increases the number of minority class samples to reduce the degree of imbalanced distribution. One of the famous over-sampling approaches is SMOTE [3]. SMOTE produces synthetic minority class samples by selecting some of the nearest minority neighbors of a minority sample which is named S , and generates new minority class samples along the lines between S and each nearest minority neighbor. SMOTE beats the random over-sampling approaches by its informed properties, and reduce the imbalanced class distribution without causing overfitting. However, SMOTE blindly generate synthetic minority class samples without considering majority class samples and may cause overgeneralization.

On the other hand, since there are much more samples of one class than the other class in the imbalanced class distribution problem, under-sampling approach is supposed to reduce the number of samples with the majority class. Assume in a training dataset, MA is the sample set with the majority class, and MI is the other set which has the minority class. Hence, an under-sampling approach is to decrease the skewed distribution of MA and MI by lowering the size of MA. Generally, the performances of over-sampling approaches are worse than that of under-sampling approaches [7].

One simple method of under-sampling is to select a subset of MA randomly and then combine them with MI as a training set, which is called *random under-sampling approach*. Several advanced researches are proposed to make the selective samples more representative. The under-sampling approach based on distance [2] uses distinct modes: the nearest, the farthest, the average nearest, and the average farthest distances between MI and MA, as four standards to select the representative samples from MA. For every minority class sample in the dataset, the first method “nearest” calculates the distances between all majority class samples and the minority class samples, and selects k majority class samples which have the smallest distances to the minority class sample. If there are n minority class samples in the dataset, the “nearest” method would finally select $k \times n$ majority class samples ($k \geq 1$). However, some samples within the selected majority class samples might duplicate.

Similar to the “nearest” method, the “farthest” method selects the majority class samples which have the farthest distances to each minority class samples. For every majority class samples in the dataset, the third method “average nearest” calculates the average distance between one majority class sample and all minority class samples. This method selects the majority class samples which have the smallest average distances. The last method “average farthest” is similar to the “average nearest” method; it selects the majority class samples which have the farthest average distances with all the minority class samples. The above under-sampling approaches based on distance in [2] spend a lot of time selecting the majority class samples in the large dataset, and they are not efficient in real applications.

In 2003, J. Zhang and I. Mani [19] presented the compared results within four informed under-sampling approaches and random under-sampling approach. The first method “*NearMiss-1*” selects the majority class samples which are close to some minority class samples. In this method, majority class samples are selected while their average distances to three closest minority class samples are the smallest. The second method “*NearMiss-2*” selects the majority class samples while their average distances to three farthest minority class samples are the smallest. The third method “*NearMiss-3*” take out a given number of the closest majority class samples for each minority class sample. Finally, the fourth method “*Most distant*” selects the majority class samples whose average distances to the three closest minority class samples are the largest. The final experimental results in [19] showed that the *NearMiss-2* method and random under-sampling method perform the best.

III. OUR APPROACH

In this section, we present our cluster-based under-sampling approach. Our approach first clusters all the training samples into some clusters. The main idea is that there are different clusters in a dataset, and each cluster seems to have distinct characteristics. If a cluster has more majority class samples and less minority class samples, it will behave like the majority class samples. On the other hand, if a cluster has more minority class samples and less majority class samples, it doesn’t hold the characteristics of the majority class samples and behaves more like the minority class samples. Therefore, our approach selects a suitable number of majority class samples from each cluster by considering the ratio of the number of majority class samples to the number of minority class samples in the cluster.

Assume that the number of samples in the class-imbalanced dataset is N , which includes majority class samples (MA) and minority class samples (MI). The size of the dataset is the number of the samples in this dataset. The size of MA is represented as $Size_{MA}$, and $Size_{MI}$ is the number of samples in MI. In the class-imbalanced dataset, $Size_{MA}$ is far larger than $Size_{MI}$. For our under-sampling method *SBC* (under-Sampling Based on Clustering), we first cluster all samples in the dataset into K clusters. In the experiments, we will study the performances for the under-sampling methods on different number of clusters. The number of majority class samples and the number of minority class samples in the i th cluster

($1 \leq i \leq K$) are $Size_{MA}^i$ and $Size_{MI}^i$, respectively. Therefore, the ratio of the number of majority class samples to the number of minority class samples in the i th cluster is $Size_{MA}^i / Size_{MI}^i$. Suppose the ratio of $Size_{MA}$ to $Size_{MI}$ in the training dataset is set to be $m:1$. The number of selected majority class samples in the i th cluster is shown in expression (1):

$$SSize_{MA}^i = (m \times Size_{MI}^i) \times \frac{Size_{MA}^i / Size_{MI}^i}{\sum_{i=1}^K Size_{MA}^i / Size_{MI}^i} \quad (1)$$

In expression (1), $m \times Size_{MI}^i$ is the total number of selected majority class samples that we suppose to have in the final training dataset. $\frac{\sum_{i=1}^K Size_{MA}^i / Size_{MI}^i}{\sum_{i=1}^K Size_{MA}^i / Size_{MI}^i}$ is the total ratio of the number of majority class samples to the number of minority class samples in all clusters. Expression (1) determines that more majority class samples would be selected in the cluster which behaves more like the majority class samples. In other words, $SSize_{MA}^i$ is larger while the i th cluster has more majority class samples and less minority class samples. If there is no minority class samples in the i th cluster, then the number of minority class samples in the i th cluster (i.e., $Size_{MI}^i$) is regarded as one, that is, we assume that there is at least one minority class sample in a cluster. After determining the number of majority class samples which are selected in the i th cluster ($1 \leq i \leq K$) by using expression (1), we randomly choose majority class samples in the i th cluster. The total number of selected majority class samples is about $m \times Size_{MI}^i$ after merging all the selected majority class samples in each cluster. Finally, we combine the whole minority class samples with the selected majority class samples to construct a new training dataset. The ratio of $Size_{MA}$ to $Size_{MI}$ is about $m:1$ in the new training dataset. Table 1 shows the steps for our cluster-based under-sampling method *SBC*.

Table 1. The structure of *SBC*

Step1.	Determine the ratio of $Size_{MA}$ to $Size_{MI}$ in the training dataset.
Step2.	Cluster all the samples in the dataset into some clusters.
Step3.	Determine the number of selected majority class samples in each cluster by using expression (1), and then randomly select the majority class samples in each cluster.
Step4.	Combine the selected majority class samples and all the minority class samples to obtain the training dataset.

For example, assume that an imbalanced class distribution dataset has totally 1100 samples. The size of MA is 1000 and the size of MI is 100. In this example, we cluster this dataset into three clusters. Table 2 shows the number of majority class

samples $Size_{MA}^i$, the number of minority class samples $Size_{MI}^i$, and the ratio of $Size_{MA}^i$ to $Size_{MI}^i$ for the i th cluster.

Table 2. Cluster descriptions

Cluster ID	Number of majority class samples	Number of minority class samples	$Size_{MA}^i / Size_{MI}^i$
1	500	10	500/10=50
2	300	50	300/50=6
3	200	40	200/40=5

Assume that the ratio of $Size_{MA}$ to $Size_{MI}$ in the training data is set to be 1:1. In other words, there are about 100 selected majority class samples and the whole 100 minority class samples in this training dataset. The number of selected majority class samples in each cluster can be calculated by expression (1). Table 3 shows the number of selected majority class samples in each cluster. We finally select the majority class samples randomly from each cluster and combine them with the minority class samples to form the new dataset.

Table 3. The number of selected majority class samples in each cluster

Cluster ID	The number of selected majority class samples
1	$1 \times 100 \times 50 / (50+6+5) = 82$
2	$1 \times 100 \times 6 / (50+6+5) = 10$
3	$1 \times 100 \times 5 / (50+6+5) = 8$

IV. EXPERIMENTAL RESULTS

In this section, we evaluate the performances for our proposed under-sampling approach on synthetic datasets. In the following, we first describe the method of generating class imbalanced datasets. And then we compare the classification accuracies of our method for minority class with the other methods by performing neural network classification algorithm [17] on synthetic datasets. Finally, the classification accuracies for minority class on real datasets by applying our proposed method and the other methods are also evaluated.

A. Generation of Synthetic Datasets

In this subsection, we present the synthetic dataset generation method to simulate the real-world dataset. This method is implemented with a user interface such that the parameters can be set for generating the synthetic dataset from the user interface, which is called *synthetic dataset generator*.

A synthetic dataset includes a set of attributes and each sample in the dataset has a set of particular attribute values. In real world, the samples in the same class should have similar attribute values and the samples in different class should have different characteristics. Even though the samples in the same

class, these samples may have different characteristics and can be clustered into some clusters. The samples in a cluster may have the similar attribute values and may belong to different classes. Besides, there may be some noises or exceptions in a dataset, that is, some samples in one class may have the similar attribute values with the samples in the other class or may be not similar to any other samples with the same class. According to the above observations, the following parameters need to be set for generating the synthetic dataset: number of samples, number of attributes and number of clusters.

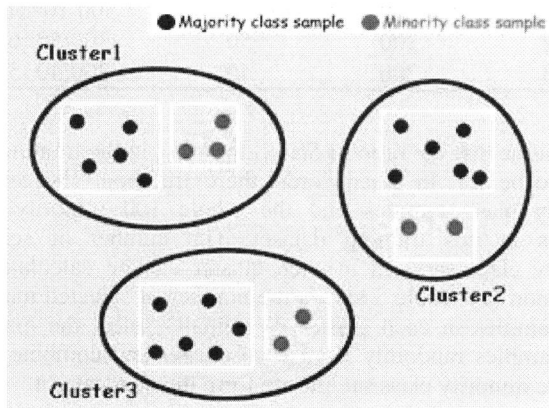


Fig1. The distribution of samples in a dataset

Because the samples in a cluster may belong to different classes, in a cluster, the samples are separated into two groups: the samples in one group are assigned a class and the samples in the other group are assigned to the other class. The attribute values for the samples are more similar to the samples in the same group, because they are in the same cluster and the same class. Fig 1 shows the distribution of samples in a dataset which has three clusters inside.

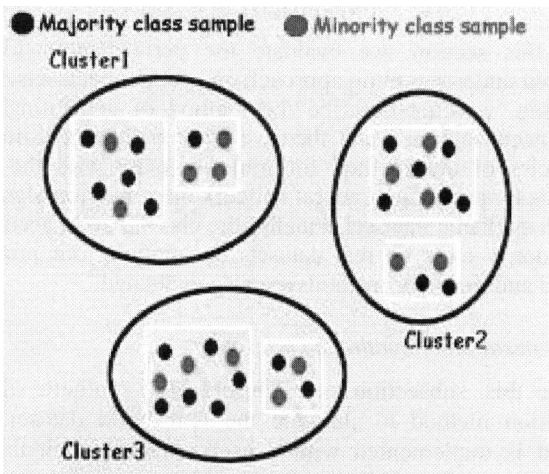


Fig 2. Example for disordered samples

In order to make the synthetic datasets more like real datasets, the noisy data are necessary. The synthetic datasets have two kinds of noisy data: disordered samples and exceptional samples. A dataset which does not have any noisy data is like the one in Fig 1. The disordered samples are

illustrated with Fig 2 in which some majority class samples (or minority class samples) lie to the area of minority class samples (or majority class samples). As for exceptional samples, they distribute irregularly in a dataset. The samples outside the clusters in Fig 3 are exceptional samples.

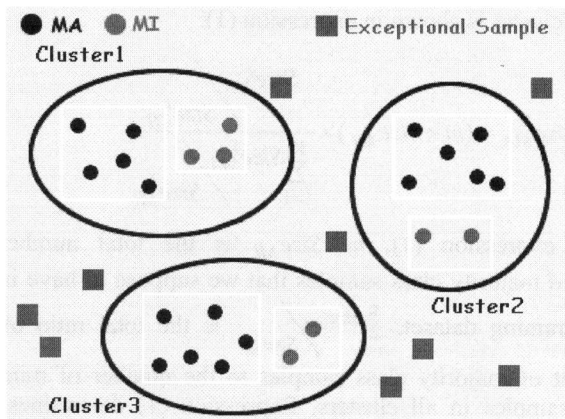


Fig 3. Example for exceptional samples

B. Evaluation Criteria

For our experiments, we use three criteria to evaluate the classification accuracy for minority class: the precision rate P , the recall rate R , and the F-measure for minority class. The precision rate for minority class is the correct-classified percentage of samples which are predicted as minority class by the classifier. The recall rate for minority class is the correct-classified percentage of all the minority class samples. Generally, for a classifier, if the precision rate is high, then the recall rate will be low, that is, the two criteria are trade-off. We cannot use one of the two criteria to evaluate the performance of a classifier. Hence, the precision rate and recall rate are combined to form another criterion F-measure, which is shown in expression (2).

$$\text{MI's F-measure} = \frac{2 \times P \times R}{P + R} \quad (2)$$

In the following, we use the three criteria discussed above to evaluate the performance of our method *SBC* by comparing our method with the other methods *AT*, *RT*, and *NearMiss-2*. The method *AT* uses all samples as the training dataset and does not select samples. *RT* is the most common-used random under-sampling method and it selects the majority class samples randomly. The last method *NearMiss-2* is proposed by J. Zhang and I. Mani [19], which has been discussed in section 2. The two methods *RT* and *NearMiss-2* have the better performance than the other proposed methods in [19]. In the following experiments, the classifiers are constructed by using the artificial neural network technique in *IBM Intelligent Miner for Data V8.1*, and the k-means clustering algorithm is used for our methods. In our experiments, the clustering algorithm would not influence the performance for our method.

C. Experimental Results on Synthetic Datasets

For each generated synthetic dataset, the number of samples is set to 10000, the number of numerical attributes and

categorical attributes are set to 5, respectively. The dataset DS_i means that the dataset potentially can be separated into *i* clusters, and our methods also cluster the dataset DS_i into *i* clusters. Moreover, a dataset DS_i with *j*% exceptional samples and *k*% disordered samples is represented as DS_iE_jD_k. If there is no disordered sample in the synthetic dataset, the dataset is represented as DS_iE_jDN.

Fig 4 shows the MI's F-measures for our method and the other methods on datasets DS4E10DN and DS4E10D20. The ratio of the number of majority class samples to the number of minority class samples is 9 to 1 in the two datasets for this experiment. In Fig 4, the method *AT* has the highest MI's F-measure in DS4E10DN because *AT* puts all the samples in the dataset into training and there is no disordered samples and just few exceptional samples in the dataset. The data distribution and characteristics can be completely represented from all the samples if there is no noise in the dataset. Hence, the classifier on DS4E10DN has the best classification accuracy when the method *AT* is applied. However, the method *AT* has to put all the samples into training, which is very time-consuming. Our method *SBC* and *RT* just need to put 20 percent of all samples into training since the ratio of $Size_{MA}$ to $Size_{MI}$ is set to be 1:1, and the MI's F-measures are above 80%. The method *AT* on dataset DS4E10D20 becomes worst and the classification accuracy is below 10%, because the dataset includes some noises, that is, 10% exceptional samples and 20% disordered samples for all the samples and all the noises are put into training. The classification accuracy for our method *SBC* and *RT* are significantly better than *AT*, since some noises can be ignored by applying *SBC* and *RT*. In this experiment, the performance of classification by using *SBC* and *RT* are better than the other methods.

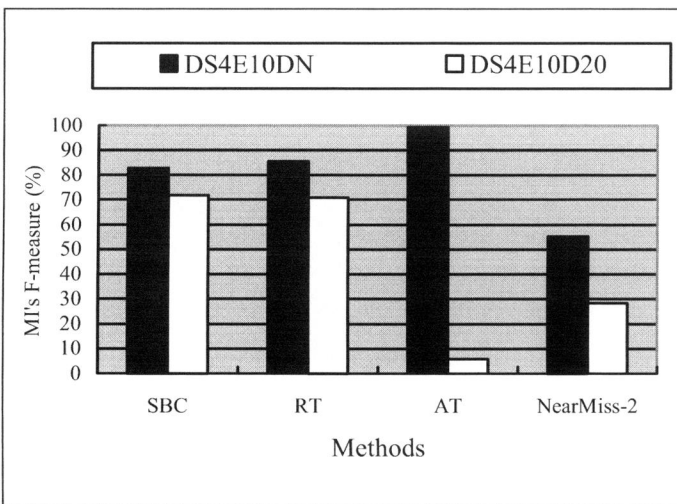


Fig 4. The effect of disordered samples

Fig 5 shows the performances of our method and the other methods on datasets DS_iE10D20, in which *i* is from 2 to 16. In these synthetic datasets, the ratio of the number of majority class samples to the number of minority class samples is 9 to 1. In Fig 5, we can see that the classification accuracy for *SBC*

and *RT* are better than other methods and our method *SBC* outperforms *RT* in most cases.

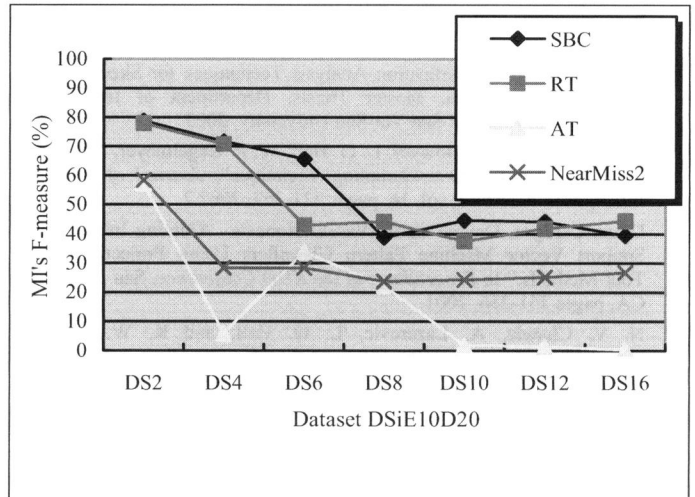


Fig 5. MI's F-measure for each method on the datasets with 10% exceptional samples and 20% disordered samples

V. CONCLUSIONS

In a classification task, the effect of imbalanced class distribution problem is often ignored. Many studies focused on improving the classification accuracy but did not consider the imbalanced class distribution problem. Hence, the classifiers which are constructed by these studies lose the ability to correctly predict the correct decision class for the minority class samples in the datasets which the number of majority class samples are much greater than the number of minority class samples. Many real applications, like rarely-seen disease investigation, credit card fraud detection, and internet intrusion detection always involve the imbalanced class distribution problem. It is hard to make right predictions on the customers or patients who that we are interested in.

In this study, we propose cluster-based under-sampling approach to solve the imbalanced class distribution problem by using backpropagation neural network. The other two under-sampling methods, Random selection and *NearMiss-2*, are used to be compared with our method in our performance studies. In the experiments, our method *SBC* has better prediction accuracy and stability than other methods. *SBC* not only has high classification accuracy on predicting the minority class samples but also has fast execution time.

ACKNOWLEDGMENT

Research on this paper was partially supported by National Science Council grant NSC 94-2213-E-130-004 and 94-2622-E-130-001-CC3.

REFERENCES

- [1] N. V. Chawla. "C4.5 and Imbalanced Datasets: Investigating the Effect of Sampling Method, probabilistic estimate, and decision tree structure." In *Proceedings of the ICML'03 Workshop on Class Imbalances*, August 2003.
- [2] Yu-Meei Chyi. "Classification Analysis Techniques for Skewed Class Distribution Problems, Master Thesis, Department of Information Management, National Sun Yat-Sen University, 2003.
- [3] N. V. Chawla, K.W. Bowyer, L. O. Hall, W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-Sampling Technique", *Journal of Artificial Intelligence Research*, vol. 16, pages 321-357, 2002.2.
- [4] Doina Caragea, Dianne Cook, Vasant Honavar. "Gaining Insights into Support Vector Machine Pattern Classifiers Using Projection-Based Tour Methods." In *Proceedings of the KDD Conference*, San Francisco, CA, pages 251-256, 2001.
- [5] N. V. Chawla, A. Lazarevic, L. O. Hall, and K. W. Bowyer. "Smoteboost: Improving Prediction of the Minority Class in Boosting." In *Proceedings of the Seventh European Conference on Principles and Practice of Knowledge Discovery in Databases*, pages 107-119, Dubrovnik, Croatia, 2003.
- [6] P. Clark and T. Niblett. "The CN2 Induction Algorithm." *Machine Learning*, 3(4):261-283, 1989.
- [7] Chris Drummond, Robert C. Holte. "C4.5, Class Imbalance, and Cost Sensitivity: Why Under-Sampling beats Over-Sampling." In *Proceedings of the ICML'03 Workshop on Learning from Imbalanced Datasets*, 2003.
- [8] C. Elkan. "The Foundations of Cost-sensitive Learning." In *Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence*, pages 973-978, 2001.
- [9] Yoav Freund, H. Sebastian Seung, Eli Shamir, Naftali Tishby. "Selective Sampling Using the Query by Committee Algorithm." *Machine Learning*, v.28 n.2-3, pages 133-168, Aug./Sept. 1997.
- [10] Rafael del-Hoyo, David Buldain, Alvaro Marco. "Supervised Classification with Associative SOM." Lecture Notes in Computer Science 2686, pp334-341. 7th *International Work-Conference on Artificial and Natural Neural Networks, IWANN 2003*.
- [11] In N. Japkowicz, editor, *Proceedings of the the AAAI'2000 Workshop on Learning from Imbalanced Data Sets*, AAAI Tech Report WS-00-05. AAAI, 2000.
- [12] N. Japkowicz. "Concept-learning in the Presence of Between-class and Within-class imbalances." In *Proceedings of the Fourteenth Conference of the Canadian Society for Computational Studies of Intelligence*, pages 67-77, 2001.
- [13] T. Jo and N. Japkowicz. "Class Imbalances versus Small Disjuncts." *SIGKDD Explorations*, 6(1):40-49, 2004.
- [14] M. Maloof. "Learning when Data Sets are Imbalanced and when Costs are Unequal and Unknown." In *Proceedings of the ICML'03 Workshop on Learning from Imbalanced Data sets*, 2003.
- [15] L. M. Manevitz and M. Yousef. "One-class SVMs for Document Classification." *Journal of Machine Learning Research*, 2:139-154, 2001.
- [16] Shlomo Argamon-Engelson, Ido Dagan. "Committee-based Sample Selection for Probabilistic Classifiers." *Journal of Artificial Intelligence Research (JAIR)*, Vol. 11, pages 335-360, 1999.
- [17] N. E. Sondak, V. K. Sondak. "Neural Networks and Artificial Intelligence." In *Proceedings of the twentieth SIGCSE technical symposium on Computer science education*. 1989.
- [18] P. Turney. "Types of Cost in Inductive Concept Learning." In *Proceedings of the ICML'2000 Workshop on Cost-Sensitive Learning*, pages 15-21, 2000.
- [19] J. Zhang and I. Mani. "kNN Approach to Unbalanced Data Distributions: A Case Study Involving Information Extraction." In *Proceedings of the ICML'2003 Workshop on Learning from Imbalanced Datasets*, 2003.