



## WPOST: Weighted Prevalence and Over Sampled Trivial of Categorical Features for Imbalanced Dataset Multi Class Learning

D. L. Rohitha  
GNITS  
Hyderabad, India

G. Malini Devi  
GNITS, Assistant Professor  
Hyderabad, India

---

**Abstract:** A dataset will be imbalanced if the classification communities won't nearly suitably represented. Commonly real-world information designs are primarily comprised of "normal" situations with just a modest proportion of "abnormal" or "interesting" designs. It is also the instance that the choice of misclassifying an extreme (interesting) instance as a frequent instance is normally greater than the cost of the treat error. Under-sampling of the more (normal) class is estimated as an effective implies of enhancing the capability of a classifier towards fraction class. In feature this, a new strategy, called weighted prevalence and over sampled trivial of specific attributes (WPOST), is revealed for efficiently dealing with imbalanced knowing issues. WPOST first acknowledges the hard-to-learn perceptive trivial class specific choices also assigns them weights which are centered on their selection from frequency class illustrations. It then generates the processed choices from the weighted worthwhile trivial class alternatives with an unsupervised learning strategy. This is achieved in such a technique that every provided selections lie within some trivial class category.

**Keywords:** ROC, AUC,

---

### I. INTRODUCTION

A dataset will be imbalanced if the classes are not inside equally symbolized. Difference on the assign of 100 to 1 is frequent in fraud identification also imbalance of upto 100,000 to 1 has get defined in other objectives (Provost & Fawcett, 2001). There have been attempts to deal with imbalanced datasets in segments like deceptive telecom management (Ezawa, Singh, & Norton, 1996), telephone calls (Fawcett & Provost, 1996), text classification (Lewis & Catlett, 1994; Mladeni'c & Grobelnik, 1999; Dumais, Platt, Heckerman, & Sahami, 1998; Lewis & Ringuette, 1994; Cohen, 1995a) also acceptance of oil spills in satellite pictures (Kubat, Holte, & Matwin, 1998).

The performance of machine reading algorithms is usually evaluated using predictive persistence. Although, this is not suggested when the information is imbalanced and/or the costs of various errors change significantly. Because an example, select the classification of pixels in mammogram images as possibly malignant (Woods, Doss, Bowyer, Solka, Priebe, & Kegelmeyer, 1993). A prevalent mammography dataset can include 98% frequent pixels as well as 2% irregular pixels. A basic standard strategy of estimating the essential class might provide a prognostic precision of 98%. Although, the type of the system needs a mainly high rate of accurate discovery in the fraction class also allows for a limited error rate in the numerous class in recommend to obtain this. Basic predictive coherence is definitely not suggested in these situations. The Receiver Operating Characteristic (ROC) curve will be a recognized strategy for detailing classifier performance over a selection of tradeoffs over true excellent as well as false excellent error levels (Swets, 1988). The Area under the Curve (AUC) will be an recognized traditional performance metric for a ROC curve (Duda, Hart, & Stork, 2001; Bradley, 1997; Lee, 2000). The ROC convex framework could also be used as a robust method of ensuring ultimately appropriate classifiers (Provost & Fawcett, 2001). If a range proceeds using a point on the convex framework, subsequently there is no assorted line using the equivalent slope going with assorted point with a significant true positive (TP) point. So, the classifier at that plan will be appropriate subsidiary any dispersal presumptions in combination with that slope.

The machine reading community has undertaken the issue of class imbalance in two methods. One will be to designate particular costs to undertake instances (Pazzani, Merz, Murphy, Ali, Hume, & Brunk, 1994; Domingos, 1999). The one other will be to re-sample the genuine dataset, possibly by oversampling the portion class and/or under-sampling the many class (Kubat & Matwin, 1997; Japkowicz, 2000; Lewis & Catlett, 1994; Ling & Li, 1998). Our technique (Chawla, Bowyer, Hall, & Kegelmeyer, 2000) combines under-sampling of the many class with a definite form of over-sampling the portion class. Experiments with various datasets as well as a Naive Bayes Classifier and the C4.5 decision tree classifier (Quinlan, 1992), Ripper (Cohen, 1995b) describe that our technique increases over other past re-sampling, maximizing loss ratio, also class priors strategies, using commonly the AUC or ROC convex framework.

### II. PERFORMANCE MEASURES

The performance of machine reading algorithms is normally evaluated by a mix-up matrix which is shown in Figure 1 (for a 2 class issue). The columns tend to be the estimated class also the rows tend to be the real class. In the distress matrix, TN will be the quantity of problem instances

	Predicted Negative	Predicted Positive
Actual Negative	TN	FP
Actual Positive	FN	TP

Figure 1: Confusion Matrix

properly classified (True Negatives), FP will be the amount of destructive instances improperly categorized as exceptional (False Positives), FN will be the amount of excellent instances improperly described as negative (False Negatives) also TP will be the amount of excellent instances properly categorized (True Positives).

Predictive reliability is the performance measure frequently associated with machine reading algorithms also is outlined as  $\text{precision} = (TP + TN) / (TP + FP + TN + FN)$ . In the view of balanced datasets as well as similar error costs, it will be moderate to utilize error rate as a performance metric. Error rate will be  $1 - \text{Accuracy}$ . In the situation of imbalanced datasets using irregular error costs, it will be most relevant to utilize the ROC curve or some other similar techniques (Ling & Li, 1998; Drummond & Holte, 2000; Provost & Fawcett, 2001; Bradley, 1997; Turney, 1996).

ROC curves could be concern of as denoting the group of best persistence limitations for equivalent costs of TP as well as FP. On an ROC curve the X-axis shows  $\%FP = FP / (TN + FP)$  and the Y-axis shows  $\%TP = TP / (TP + FN)$ . The appropriate point on the ROC curve could be (0,100), that is every assenting instance tend to be classified efficiently also no negative designs are misclassified as constructive. One strategy an ROC curve could be obtained is by altering the balance of expertise remedies for every class for the training set. The line  $y = x$  represents the situation of arbitrarily estimating the class. Area inside the ROC Curve (AUC) will be an efficient metric for classifier performance because it is definite of the persistence standard selected also prior possibilities. The AUC evaluation can ascertain a visibility connection perhaps classifiers. If the ROC curves tend to be intersecting, the overall AUC will be a typical evaluation among products (Lee, 2000). While, for certain specific cost as well as class distributions, the classifier acquiring optimum AUC can in actuality become suboptimal. Hence, we even determine the ROC convex hulls, as the points lying from the ROC convex framework tend to be extremely optimum (Provost, Fawcett, & Kohavi, 1998; Provost & Fawcett, 2001).

### III. WPOST: WEIGHTED PREVALENCE AND OVER SAMPLED TRIVIAL

#### WPOST

We provide an over-sampling technique wherein the portion class will be over-sampled by generating “synthetic” designs rather than by over-sampling using substitution. This technique is established by a strategy that revealed efficient in written character recognition (Ha & Bunke, 1997). They generated additional training information by performing specific operations on considerable data. In their situations, procedures such as sequence as well as skew were normal techniques to perturb the training data. We generate synthetic designs in a reduced application-specific method, by executing in “feature space” rather than “data space”. The portion class will be over-sampled by obtaining every fraction class test also evaluating synthetic designs over the line segments connecting any/all of the k portion class nearest neighbors. Considering the quantity of over-sampling required, neighbors with the k nearest neighbors tend to be randomly selected. Our performance currently utilizes five nearest neighbors. For example, if the quantity of over-sampling needed is 200%, just two neighbors from the five nearest neighbors tend to be chosen as well as one sample will be provided in the position of all.

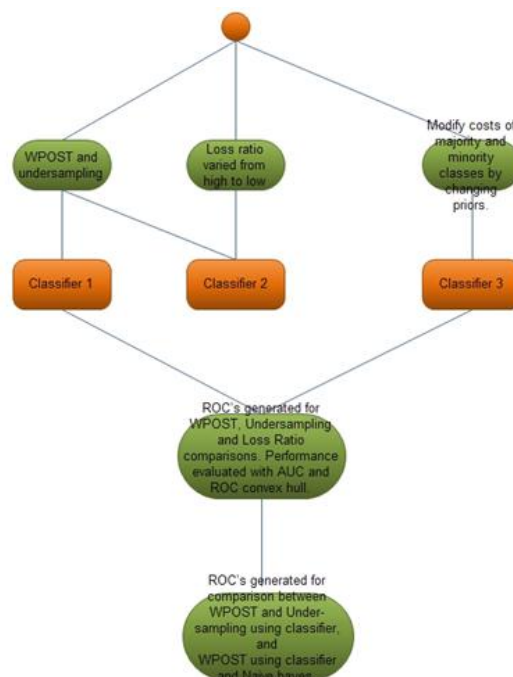


Fig2: Proposed Architecture for Weighted prevalence and over sampled trivial (WPOST)

Artificial products are developed in the suitable strategy: choose the variation among the component vector (sample) here focus also its nearest neighbor. Enhance this transform by a distinctive number around 0 and 1, also add it to the component vector under consider. This generates the variety of a decisive point collectively the line segment around two separate features. This technique effectively forces the persistence region of the portion class to get more consistent.

Algorithm WPOST, in the appropriate page, will the pseudo-code to WPOST. The sum of over-sampling will be a amount of the program, also a set of ROC curves may be provided for numerous populations as well as ROC assessment performed. The synthetic instances trigger the classifier towards generate larger as well as fewer particular persistence regions, rather than restrained also additional certain regions. Most usual regions are actually determined for the fraction class choices instead of those obtaining subsumed by many class selections around them. The affect is that persistence trees generalize adept. The studies were carried out on the mammography dataset. There comprise 10923 cases within the majority class also 260 cases in the minority class originally. We've about 9831 cases in the majority class also 233 cases in the minority class towards training set used at 10-fold cross-validation. The portion class got over-sampled at 100%, 200%, 300%, 400% as well as 500% of its unique size. The chart reveal that the tree dimensions for portion over-sampling using substitution at greater degrees of replication are a lot superior to those for WPOST, also the minority class recognition from the minority over-sampling using substitution strategy at greater degrees of replication isn't as effective as WPOST.

**Algorithm WPOST** ( $T, N, k$ )

**Input:** Amount of minority class samples  $T$ ; Number of WPOST  $N\%$ ; Amount of nearest neighbors  $k$

**Output:**  $(N/100)*T$  synthetic minority class samples

1. (\*If  $N$  is less than 100%, randomize the minority class samples as only a random percent of them will be SMOTEd.\*)
2. **if**  $N < 100$
3. **then** Randomize the  $T$  minority class samples
4.  $T = (N/100)*T$
5.  $N = 100$
6. **endif**
7.  $N = (\text{int})(N/100)$  (\*The amount of SMOTE is assumed to be in integral multiples of 100.\*)
8.  $k =$  Number of nearest neighbors
9.  $\text{numattrs} =$  Number of attributes
10.  $\text{Sample} [ \ ] [ \ ]$ : array for original minority class samples
11.  $\text{newindex}$ : keeps a count of number of synthetic samples generated, initialized to 0
12.  $\text{Synthetic} [ \ ] [ \ ]$ : array for synthetic samples  
(\*Compute  $k$  nearest neighbors for each minority class sample only.\*)
13. **for**  $i \leftarrow 1$  **to**  $T$
14. Compute  $k$  nearest neighbors for  $i$ , and save the indices in the  $\text{marray}$
15.  $\text{Populate}(N, i, \text{marray})$
16. **endfor**  
 $\text{Populate}(N, i, \text{marray})$  (\*Function to generate the synthetic samples.\*)
17. **while**  $N \neq 0$
18. Choose a random number between 1 and  $k$ , call it  $mn$ . This step chooses one of the  $k$  nearest neighbors of  $i$ .
19. **for**  $\text{attr} \leftarrow 1$  **to**  $\text{numattrs}$
20. Compute:  $\text{dif} = \text{Sample}[\text{marray}[mn]][\text{attr}] - \text{Sample}[i][\text{attr}]$
21. Compute:  $\text{gap} =$  random number between 0 and 1
22.  $\text{Synthetic}[\text{newindex}][\text{attr}] = \text{Sample}[i][\text{attr}] + \text{gap} * \text{dif}$
23. **endfor**
24.  $\text{newindex}++$
25.  $N = N - 1$
26. **endwhile**
27. **return** (\*End of Populate.\*)  
End of Pseudo-Code

#### IV. UNDER-SAMPLING AND WPOST COMBINATION

The majority class will be under-sampled by arbitrarily elimination of selections from the majority class group till the fraction class receives most specific amount of the majority class. This aspects the learner to undertake altering

degrees of under-sampling as well as at greater degrees of under-sampling the portion class offers a significant situation in the training set. In detailing our studies, our words are such as that if we under-sample the majority class at 200%, it could indicate that the personalized dataset can contain twice as many components from the portion class as within the majority class; which, if the portion class had 50 choices also the most class had 200 choices as well as we under-sample most at 200%, the most class may end up with 25 samples. By utilizing an incorporating of under-sampling as well as over-sampling, the initial tendency of the learner about the negative (majority) class will be inverted in the assistance of the positive (minority) class. Classifiers tend to be determined on the dataset perturbed with “SMOTING” the portion class as well as under-sampling the majority class.

## V. EXPERIMENTAL RESULTS

A ROC curve for WPOST will be provided by utilizing C4.5 or Ripper to generate a classifier for all of a series of altered training datasets. A provided ROC curve will be provided by initial over-sampling the fraction class to a certain degree thereafter under-sampling the most class at maximizing degrees to produce the consecutive points in the curve. The quantity of under-sampling will be equivalent to simply under-sampling. Hence, every equivalent point on every ROC curve with a dataset signifies the equivalent amount of most class samples. Various ROC curves tend to be provided by initiating with various levels of fraction over-sampling. ROC curves were even produced by changing the loss ratio within Ripper from 0.9 to 0.001 also by changing the priors of the fraction class from the initial delivery around 50 times the most class for a Naive Bayes Classifier.

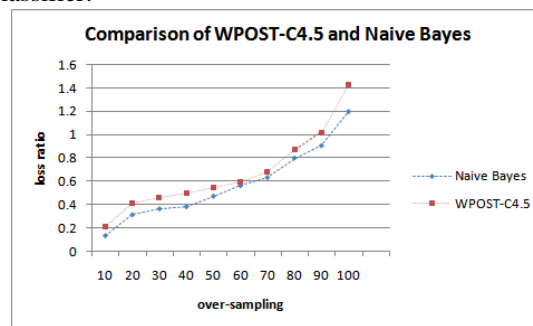


Figure 3: Phoneme. Comparison of WPOST-C4.5, Under-C4.5, and Naive Bayes. WPOST C4.5 dominates over Naive Bayes and Under-C4.5 in the ROC space. WPOST C4.5 classifiers are potentially optimal classifiers.

## VI. CONCLUSION

There are several topics being considered additional in this range of evaluation. Automated adaptive collection of the amount of nearest neighbors could be efficient. Assorted techniques for producing the synthetic neighbors may be capable of elevate the performance. Furthermore, choosing closest neighbors with a focus on instances that are inappropriately classified may elevate performance. A fraction class test might perhaps have a most class test as its nearest neighbor rather than a fraction class sample. This crowding may possibly provide to the redrawing of the persistence areas for the portion class. In order to these concerns, the suitable subsections consult two obtainable extensions of WPOST, also a system of WPOST to information retrieval.

## REFERENCES

- [1] Blake, C., & Merz, C. (1998). UCI Repository of Machine Learning Databases <http://www.ics.uci.edu/~mllearn/~MLRepository.html>. Department of Information and Computer Sciences, University of California, Irvine.
- [2] Bradley, A. P. (1997). The Use of the Area Under the ROC Curve in the Evaluation of Machine Learning Algorithms. *Pattern Recognition*, 30(6), 1145–1159.
- [3] Chawla, N., Bowyer, K., Hall, L., & Kegelmeyer, P. (2000). SMOTE: Synthetic Minority Over-sampling TEchnique. In *International Conference of Knowledge Based Computer Systems*, pp. 46–57. National Center for Software Technology, Mumbai, India, Allied Press.
- [4] Chawla, N., & Hall, L. (1999). Modifying MUSTAFA to capture salient data. Tech. rep. ISL-99-01, University of South Florida, Computer Science and Eng. Dept.
- [5] Cohen, W. (1995a). Learning to Classify English Text with ILP Methods. In *Proceedings of the 5th International Workshop on Inductive Logic Programming*, pp. 3–24. Department of Computer Science, Katholieke Universiteit Leuven.
- [6] Cohen, W. W. (1995b). Fast Effective Rule Induction. In *Proc. 12th International Conference on Machine Learning*, pp. 115–123 Lake Tahoe, CA. Morgan Kaufmann.
- [7] Cohen, W. W., & Singer, Y. (1996). Context-sensitive Learning Methods for Text Categorization. In Frei, H.-P., Harman, D., Schöuble, P., & Wilkinson, R. (Eds.), *Proceedings of SIGIR-96, 19th ACM International Conference on Research and Development in Information Retrieval*, pp. 307–315 Zürich, CH. ACM Press, New York, US.
- [8] Cost, S., & Salzberg, S. (1993). A Weighted Nearest Neighbor Algorithm for Learning with Symbolic Features. *Machine Learning*, 10(1), 57–78.

- [9] DeRouin, E., Brown, J., Fausett, L., & Schneider, M. (1991). Neural Network Training on Unequally Represented Classes. In *Intelligent Engineering Systems Through Artificial Neural Networks*, pp. 135–141 New York. ASME Press.
- [10] Domingos, P. (1999). Metacost: A General Method for Making Classifiers Cost-sensitive. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 155–164 San Diego, CA. ACM Press.
- [11] Drummond, C., & Holte, R. (2000). Explicitly Representing Expected Cost: An Alternative to ROC Representation. In *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 198–207 Boston. ACM.
- [12] Duda, R., Hart, P., & Stork, D. (2001). *Pattern Classification*. Wiley-Interscience.
- [13] Hall, L., Mohny, B., & Kier, L. (1991). The Electrotopological State: Structure Information at the Atomic Level for Molecular Graphs. *Journal of Chemical Information and Computer Science*, 31(76).
- [14] Japkowicz, N. (2000). The Class Imbalance Problem: Significance and Strategies. In *Proceedings of the 2000 International Conference on Artificial Intelligence (IC-AI'2000): Special Track on Inductive Learning Las Vegas, Nevada*.
- [15] Kubat, M., Holte, R., & Matwin, S. (1998). Machine Learning for the Detection of Oil Spills in Satellite Radar Images. *Machine Learning*, 30, 195–215.
- [16] Kubat, M., & Matwin, S. (1997). Addressing the Curse of Imbalanced Training Sets: One Sided Selection. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pp. 179–186 Nashville, Tennessee. Morgan Kaufmann.
- [17] Lee, S. (2000). Noisy Replication in Skewed Binary Classification. *Computational Statistics and Data Analysis*, 34.
- [18] Lewis, D., & Catlett, J. (1994). Heterogeneous Uncertainty Sampling for Supervised Learning. In *Proceedings of the Eleventh International Conference of Machine Learning*, pp. 148–156 San Francisco, CA. Morgan Kaufmann.
- [19] Quinlan, J. (1992). *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo, CA.
- [20] Solberg, A., & Solberg, R. (1996). A Large-Scale Evaluation of Features for Automatic Detection of Oil Spills in ERS SAR Images. In *International Geoscience and Remote Sensing Symposium*, pp. 1484–1486 Lincoln, NE.