

Improving Labeling Quality using Positive Label Frequency Threshold Algorithm

M.Kiruthiga¹, P.Sangeetha²

¹PG Scholar, ²Assistant Professor,

^{1,2}Department of Computer Science and Engineering

^{1,2}P. A. College of Engineering and Technology, Pollachi, Coimbatore, India
kiruthigamanickam31@gmail.com¹, sangeethasiva79@gmail.com²

Abstract: Label is a prominent issue in the classification area along with several potential negative sequences. For example, the predicted accuracy may reduce, but the complexity of inferred models and the number of necessary training samples may rise. Online outsourcing systems, such as Amazon's Mechanical Turk, allow labelers to label the same objects but still lack in their quality. Mostly noisy labels have multiple labels for same examples. Thus, an agnostic algorithm Positive LABEL frequency Threshold (PLAT) is projected to handle the issue of imbalanced noisy labeling. The main objective is to generate the training dataset and integrated labels of examples. This method is used to solve the issue of minority sample and also able to deal with imbalanced multiple noisy labeling. The PLAT is applied to the imbalanced dataset collected from Amazon Mechanical Turk and the experiment results represents that the PLAT is efficient than other methods.

Index Terms –repeated labeling, majority voting, imbalanced labeling

I. INTRODUCTION

The online crowd sourcing systems such as Rent-A-Coder and Amazon Mechanical Turk, is to acquire required services, generating ideas or content from a large group of people, allowing number of labelers to label the object and facilitates to perform non-expert labeling inexpensively and simply. Thus online crowd sourcing are gainful while comparing to traditional expert labeling methods. The cheap labels are noisy due to missing of the expertise, preference and enthusiasm. Considering repeated labeling is determining multiple labels for all data points [11]. preceding research describes repeated labeling strategies can improve the labeling quality by integrating the repeated labels using Majority Voting (MV) integration strategy. For example, considering a multiple noisy label set {+, -, +, -, +} and applying the MV, as a result final label “+” is assigned to this example since “+” obtains the highest voting. A preceding scenario strategy of using multiple noisy labels Majority Voting (MV) classify the class based on the highest number of prediction or voting received. It integrates the labels and assumes that all data points are distributed uniformly and completed that the quality of labels are higher. However, the mislabeling is generally not distributed uniformly. Taking binary classification labelers provide higher probability for one and other with lesser probability [10]. For example, mostly labeling on minority examples is error-prone and it is not unusual. In this scenario, the algorithm handles minority as the positive class. While the labels are imbalanced, the number of negative labels obtained is far more than that of positive labels. The negative examples outnumber positive ones when MV is applied and the training set contain no positive examples.

We introduced an agnostic algorithm PLAT to use skewed noisy labels to stimulate an integrated label for each example. Imbalanced noisy labeling make the minority examples (assuming positive cases) in the training set occur rarely and this problem is solved by this algorithm. The organization of the paper is as follows. In section 2, the related works are reviewed. In Section 3, the estimation of accuracy are analyzed. Section 4 describes the working of an agnostic algorithm. In section 5, we compare the performance of our algorithm with other method. Section 6 provides the conclusion and future work.

II. RELATED WORK

A imbalanced datasets is learned based on a combination of the SMOTE algorithm and the boosting procedure to improve the overall F-values and to get better prediction performance on the minority class [2]. He et al. evaluated the learning performance under the imbalanced learning scenario by provided a review on the state-of-the-art technologies, and the current assessment metrics [5]. Donmez presented Interval Estimate (IE) Threshold to pick the experts with the highest accuracy for labels [3]. Kumar defines the supervised learning methods where unsupervised counter-parts are outperformed frequently since the learner are provided with more information can permit to learn a desired pattern effectively [7]. Smyth et al. described the remote sensing applications, to train the pattern recognition algorithms to detect objects of interest using ground-truth data as the basis [13]. [6] Kajino et al. proposed a convex optimization formulation for learning from crowds .The personal models for each individual crowd workers are introduced for estimating without the true labels. Strapparava et al. presents the Affective Text task and it is focused on the labeling of emotions and valence classification in news headlines, and is intended as an exploration of the connection between emotions and lexical semantics [15]. Lo et al., describe the Cost-Sensitive learning problems [9]. It is based on audio tag annotation task and by treating the tag count as costs the cost sensitive classification issues are solved. Our work is different and the examples are given the higher priority. Two classes are treated equally in our work.

III. ACCURACY ESTIMATION

The proportion of both true positives and true negatives among the total number of cases is described as accuracy and it is examined . The minority class is used as positive class and majority class as negative class, we calculate the accuracy using following equation (1),

$$\text{Accuracy} = \frac{\text{True positive} + \text{True negative}}{\text{Total number of true cases}} \quad (1)$$

The true positive (TP) is the number of correctly labeled items that belong to the positive class. The true negative (TN) is the number of correctly labeled items that belong to the negative class. The false positive (FP) is the number of items incorrectly labeled as belonging to the positive class. The false negative (FN) is the number of items incorrectly labeled as belonging to the negative class. The accuracy is the evaluation of classifier on a set of test data. Based on the number of instances in the test data, the correct

classifiers prediction is found. The 100% accuracy defines that the provided value and the measured values are accurately the same.

A. IMBALANCED LABELING IMPACT ON MV

A data set containing a proportion tp of true positive examples and tn of true negative examples is considered, the class distribution is balanced if $tp \leq 0.5$. A variable V is defined to control the mislabeling percentage on the positive data points. It reflects the imbalanced labeling level, the higher level of imbalance. The labeling quality can be integrated on positive examples Pp , and Pn on negative examples if the labeling quality is same for all labelers, then $Pp = (tp + Vp - V)/d$ and $Pn = (p + V - Vp - tp)/(tn)$ are calculated. When applying the majority voting, we can use Bernoulli model to calculate the integrated quality q of multiple noisy labels by using. Then α which is the ratio of the labeled number of positive examples (Pos) and negative examples (Neg) are evaluated as follows,

$$\alpha = \text{Pos/Neg} = [tpq_p + (1-d)(1-qn)] / [(tn)qn + tp(1-qp)] \quad (2)$$

For example, if $d=0.5$ the class distribution is balanced and if $0.5 < p < 1$, increasing number of labels decreases α value by applying MV . Thus the accuracy of learning model will eventually decrease when α is reduced and the number of positive examples in the final training set will also decline. It gives rise imbalanced noisy labeling and also results in low quality labeling. Thus if the distribution of class is imbalance, then the outcome will be worse. Thus MV is easy to understand but for imbalanced multiple noisy labeling, the MV does not work [17] at all. Certain sampling techniques [4], [8] may also be used but the limitation over that method is the important information also gets eliminated.

IV. PLAT ALGORITHM

The threshold algorithm is to check and create an effective label for multiple noisy label dataset. In mushroom dataset [1], considering a specific sample $s_i = \langle x_i, y_i \rangle$ and it associates a multiple noisy label set that enclose $L_{pos}^{(i)}$ positive labels and $L_{neg}^{(i)}$ negative labels. Using it the frequency of positive and negative labels are determined using equation (3),

$$\text{Freq}_p = \frac{L_{pos}}{L_{pos} + L_{neg}} = 1 - \text{Freq}_n \quad (3)$$

To obtain the efficient result, we introduced the technique Positive Label Frequency threshold (PLAT) algorithm to process the noisy dataset more effectively. The sample set is considered as input that contains the examples with multiple noisy label set. Finally, as a output the positive and negative are listed.

Algorithm

1. For each $i \in \text{Sample_set}$ do
2. Calculate F_i and insert it into frequency_table
3. Initialize final labels of samples to be negative

4. Sort (frequency_table) in ascending order of F
5. $N_0 := \text{size of (sample set)}$
6. $N_{L1} := N_{R1} := 0$
7. $P = \text{EstimateThresholdPosition}(\text{frequency_table}, N_0, N_{L1}, N_{R1})$
8. $P_{0_{max}} = (N_{L1} - N_{R1}) * N_{R1} / (N_{L1} + N_{R1}) + N_{R1}$
9. $L = \text{size of}(\text{frequency_table}) - 1$
10. $Pos = 0$
11. While $L > P$ do
12. Category (F_L) = pos
13. $Pos = Pos + \text{sizeof}(\text{items}(F_L))$
14. $F_m = (F_0 + F_t) * \theta$
15. $L = P$
16. While
17. $F_L > F_m \& Pos + \text{sizeof}(\text{items}(F_L)) < P_{0_{max}}$ do
18. Category (F_L) = pos
19. $Pos = Pos + \text{sizeof}(\text{items}(F_L))$
20. For $i = 0$ to $\text{size of}(\text{frequency_table}) - 1$ do
21. Insert $\text{items}(F_i)$ into list_p or list_n according to category (F_L) value
22. return list_p and list_n.

Initially we have to split the given frequency table into multiple range intervals. We can directly classify the samples whose values are greater than the specified threshold T value as positive samples. For interval with the Freq_p values less than and equal to the threshold T, we compute the middle value of Freq_p of this interval F_m . We consider that the data points whose Freq_p values are greater than f_m and close to threshold T have high probability to be positive. The remaining data points have high probability to be negative. The algorithm shows that the category of the data points is found by the final proportion of the positive and negative cases. Finally the algorithm return the positive and negative lists and the accuracy is calculated. Thus, the algorithm solve the problem and improve the label quality [12].

A. ESTIMATE THRESHOLD POSITION ALGORITHM

The sorted frequency_table, N_0 as input, the position P in sorted frequency_table whose value is treated as threshold T is evaluated.

1. Add position 0 to max_set
2. For $i = 1$ to $\text{size of}(\text{frequency_table}) - 2$ do
3. $a_0 = \text{size of}(\text{items}(F_i)) - \text{sizeof}(\text{items}(F_i - 1))$
4. $b_0 = \text{sizeof}(\text{items}(F_i + 1)) - \text{sizeof}(\text{items}(F_i))$
5. if $a_0 \geq 0 \& b_0 \leq 0$ $\text{diff}(F_i, F_{\text{min_set}(\text{last})})$,
6. then add i into max_set
7. if $a_0 \leq 0 \& b_0 \geq 0$ $\text{diff}(F_{\text{min_set}(\text{last})}, F_i)$,
8. then add i into minima_set
9. $P_0 = \text{argmax}_j \{ \text{sizeof}(\text{items}(F_j)) \mid F_j < 0.5, j \in \text{max_set} \}$
10. $P_1 = \text{argmax}_k \{ \text{sizeof}(\text{items}(f_j)) \mid f_j < 0.5, k \in \text{max_set} \}$
11. If P_0 & P_1 are not found then valley = $\text{argmin}_l \{ \text{sizeof}(\text{items}(F_l)) \mid F_{P_0} < F_l < F_{P_1}, l \in \text{min_set} \}$
12. If valley found then $P = \text{valley}$ else $P = P_0$
13. $N_{L1} = \sum_{i=0}^P \text{sizeof}(\text{items}(F_i))$

14. While $N_{L1} < N0/2$ do
15. $P=P+1$;
16. $N_{L1} = N_{L1} + \text{sizeof}(\text{items}(Fp))$
17. $N_{R1}=N0-N_{L1}$
18. Return P, N_{L1} and N_{R1}

This algorithm describes that the probability of positive and negative sample computation and return the position for each example and then it is used in the PLAT algorithm to list the positive and negative cases.

V. EXPERIMENTS

The performance of PLAT algorithm is estimated on conducting experiment on mushroom dataset listed in Table 1. The mushroom dataset includes hypothetical samples corresponding to 23 species of gilled mushrooms in the agaricus and lepiota family. Each species is identified as edible, poisonous.

TABLE 1 Dataset Used in Experiment

Dataset	Mushroom
Attributes	23
Examples	8124
Positive	3916
Negative	4208

The mushroom dataset is shown in Fig 5.1. The PLAT algorithm is based on the distribution of positive and negative labels and the accuracy is calculated for it based on the labeling. Each non-numeric attribute is converted to numeric values and the missing attribute is assigned to zero. Then, the positive and negative labels are assigned to the each tuples. Using majority voting, the accuracy is evaluated and is shown in Fig 5.2. The position is estimated to each elements as shown in Figure 5.3. Then the accuracy is evaluated for PLAT algorithm and is shown in Figure 5.6.

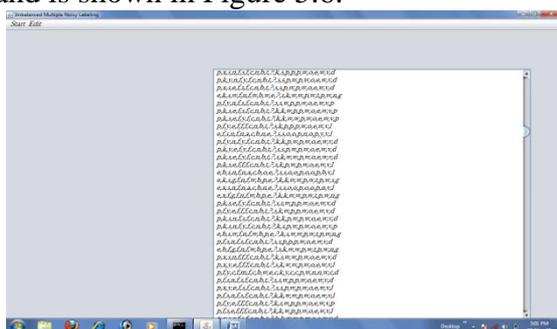


Fig 5.1 Mushroom dataset

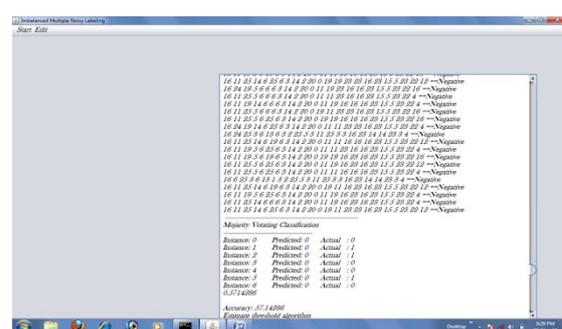


Figure 5.2 Accuracy of Majority Voting

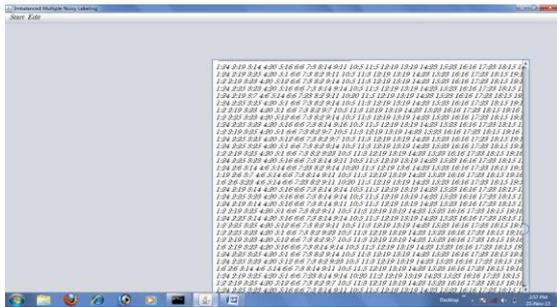


Figure 5.3 Position Estimation using Estimate Threshold Position Algorithm.

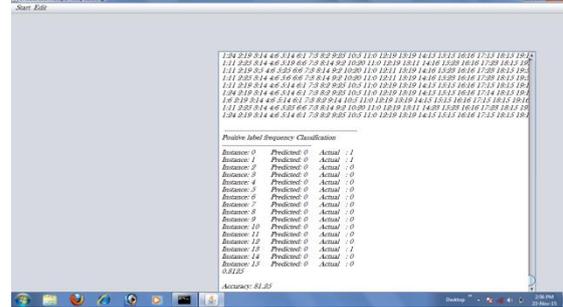


Figure 5.4 Accuracy estimation for PLAT Algorithm

Finally, the PLAT algorithm is compared with majority voting method. Under imbalanced class distribution, the performances of both methods are evaluated. Thus PLAT algorithm is agnostic and it produces the highest accuracy value when comparing with the majority voting is shown in Table 2.

TABLE 2 Performance comparison on mushroom dataset

METHOD	ACCURACY
MV Method	64.5
PLAT Algorithm	81.5

VI. CONCLUSION

The PLAT algorithm effectively performs the imbalanced labeling and thus it does not require any knowledge of labelers labeling quality and thus it can be used for both balanced and imbalanced labeling. The experimental result shows that it performs well and in future the cost-sensitive learning can be studied to evaluate the misclassification cost.

REFERENCE

- [1] C. L. Black and C. J. Merz. UCI repository of machine learning database [Online]. Available: <http://archive.ics.uci.edu/ml/>, 1998.
- [2] N. V. Chawla, L. O. Hall, K. W. Bowyer, and W. P. Kegelmeyer, "SMOTE: Synthetic minority oversampling technique," J. Artif. Intell. Res., vol. 16, pp. 321–357, 2002.
- [3] P. Donmez, J. G. Carbonell, and J. Schneider, "Efficiently learning the accuracy of labeling sources for selective sampling," in Proc.15th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2009, pp. 259–268.
- [4] A. Estabrooks, T. Jo, and N. Japkowicz, "A multiple resampling method for learning from imbalanced data sets," Comput. Intell., vol. 20, no. 1, pp. 18–36, 2004.

- [5] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 9, pp. 1263–1284, Sep. 2009.
- [6] H. Kajino, Y. Tsuboi, and H. Kashima, "A convex formulation for learning from crowds," in *Proc. 26th AAAI Conf. Artif. Intell.*, 2012, pp. 73–79.
- [7] A. Kumar and M. Lease, "Modeling annotator accuracies for supervised learning," in *Proc. 4th ACM WSDM Workshop Crowd sourcing Search Data Mining*, 2011, pp. 19–22.
- [8] X. Y. Liu, J. Wu, and Z. H. Zhou, "Exploratory under sampling for class imbalance learning," in *Proc. IEEE 6th Int. Conf. Data Mining*, 2006, pp. 965–969.
- [9] H. Y. Lo, J. C. Wang, H. M., Wang, and S. D., Lin, "Cost-sensitive multi-label learning for audio tag annotation and retrieval," *IEEE Trans. Multimedia*, vol. 13, no. 3, pp. 518–529, Jun. 2011.
- [10] C. Parker, "On measuring the performance of binary classifiers," *Knowl. Inform. Syst.*, vol. 35, no. 1, pp. 131–152, 2013.
- [11] V. S. Sheng, "Simple multiple noisy label utilization strategies," in *Proc. IEEE 11th Int. Conf. Data Mining*, 2011, pp. 635–644.
- [12] V. S. Sheng, F. Provost, and P. Ipeirotis, "Get another label? Improving data quality and data mining using multiple, noisy labeler," in *Proc. 14th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2008, pp. 614–662.
- [13] P. Smyth, M. C. Burl, U. M. Fayyad, P. Perona, and P. Baldi, "Inferring ground truth from subjective labeling of venus images," *Adv. Neural Inform. Process. Syst.*, vol. 8, pp. 1085–1092, 1995.
- [14] R. Snow, B. O'Connor, D. Jurafsky, and A. Ng, "Cheap and fast— But is it good?" in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2008, pp. 254–263.
- [15] C. Strapparava and R. Mihalcea, "SemEval-2007 Task 14: Affective text," in *Proc. 4th Int. Workshop Semantic Eval.*, 2007, pp. 70–74.
- [16] P. Welinder and P. Perona, "Online crowdsourcing: Rating annotators and obtaining cost-effective labels," in *Proc. Workshop Adv. Comput. Vis. Humans Loop*, 2010, pp. 25–32.
- [17] J. Whitehill, P. Ruvolo, T. Wu, J. Bergsma, and J. Movellan, "Whose vote should count more: Optimal integration of labels from labelers of unknown expertise," in *Proc. Adv. Neural Info. Process. Syst.* 22, 2009, pp. 2035–2043.
- [18] J. Zhang, X. Wu, and Victor S. Sheng, "Imbalanced Multiple Noisy Labeling", vol 27, feb 2015.