

## DW-ML-kNN: A Dual Weighted Multi-label kNN Algorithm

Duoqian Miao, Zhifei Zhang<sup>\*</sup>, Zhihua Wei, Chuanyan Wang

*Department of Computer Science and Technology, Tongji University, No. 4800, Cao'an Highway  
Shanghai, 201804, China  
{dqmiao, zhihua\_wei}@tongji.edu.cn, zhifei.zzh@gmail.com, chunyanfriend@163.com*

Multi-label learning is a continually developing field in machine learning and has wide applications in many domains. ML-kNN is a simple but effective method to learn from multi-label data. But when data are imbalanced, its performance is not encouraging. Therefore, we propose a dual weighted ML-kNN (for short, DW-ML-kNN) algorithm, which utilizes distances of neighbors as weights and considers the impact of neighbors without one certain label. The experiment results on three multi-label datasets show that the new algorithm achieves better performance than ML-kNN.

*Keywords:* Multi-label learning; k nearest neighbors; distance weight.

### 1. Introduction

Multi-label data are popular in our daily life. For example, an image may be classified into semantic classes *sunset* and *beach* simultaneously; a news report is about *politics*, as well as *military*. Recently, the issue of learning from these data, which is called multi-label learning, has drawn significant attention of researchers. Multi-label learning is widely used in many domains, such as text classification<sup>1,2,3</sup>, semantic annotation of videos<sup>4</sup>, image annotation<sup>5,6</sup> and functional genomics<sup>7,8</sup>. There are two major tasks in multi-label learning: multi-label classification and label ranking<sup>9</sup>. The methods of solving this problem are grouped into two categories: problem transformation and algorithm adaptation<sup>10</sup>.

Problem transformation is to transform the multi-label learning problem into one or more single-label learning problem, for which traditional learning algorithms are used. Ranking by pairwise comparison proposed by Hüllermeier et al.<sup>11</sup> constructs multiple binary classifiers, and determines the label ranking of a new instance by voting. Calibrated label ranking<sup>12</sup> extends it by introducing an additional virtual label, which assigns the labels ranking before the virtual label to unseen example. Zhang et al.<sup>13</sup> gives an algorithm called INSDIF, in which each prototype for each label is computed with the all instances belonging to this label as the training set. This group of methods poses an important complexity problem, especially for large value of the number of labels. PPT<sup>14</sup> and RAKE<sup>15</sup> methods solve the problem to some extent.

---

<sup>\*</sup> Corresponding Author

The comparison results on Emotions dataset are described in Fig. 4. The meaning of the coordinates is the same with that in Fig. 2. The performance of DW-ML-kNN on Precision and Recall is better than ML-kNN, but not obviously on other criteria, especially worse on Coverage. The reason is similar to that of Scene dataset.

The overall result is obtained from the above three results. Table 3 reports the average performance of the two algorithms across 11 values of the parameter  $K$  for each dataset. The best result on each criterion is shown with bold typeface. The last line means the number of criteria on which each algorithm achieves the best result.

Table 3. Overall Comparison Results of Two Algorithms

Criterion	Yeast		Scene		Emotions	
	ML-kNN	DW-ML-kNN	ML-kNN	DW-ML-kNN	ML-kNN	DW-ML-kNN
Hamming loss	<b>0.1973</b>	0.2028	<b>0.0978</b>	0.1013	<b>0.2900</b>	0.2981
Precision	<b>0.7273</b>	0.6880	0.6583	<b>0.7106</b>	0.4928	<b>0.5468</b>
Recall	0.5653	<b>0.6275</b>	0.6471	<b>0.7625</b>	0.3510	<b>0.4781</b>
One-error	0.2405	<b>0.2398</b>	0.2481	<b>0.2462</b>	<b>0.3896</b>	0.4012
Coverage	6.4015	<b>6.3956</b>	<b>0.5804</b>	0.5828	2.4980	<b>2.4888</b>
Average precision	0.7576	<b>0.7595</b>	0.8472	<b>0.8490</b>	0.6898	<b>0.6994</b>
Win(s)	2(6)	4(6)	2(6)	4(6)	2(6)	4(6)

The results show that DW-ML-kNN is better than ML-kNN on more than half of the six criteria on three benchmark datasets. Average precision and Recall are always improved, while Hamming loss is still worse. The performance on other three criteria is unstable, but close. In a word, the dual weighted ML-kNN (DW-ML-kNN) algorithm achieves promising and better performance.

## 5. Conclusion

To solve the problem of imbalanced data, we propose a new algorithm-dual weighted multi-label k nearest neighbors, named DW-ML-kNN. The experiment results show that our algorithm is better than ML-kNN. We will carry out research on the criterion of Hamming loss, and try to improve its performance further. What's more, the tuning factor in our algorithm is also a further research topic.

## Acknowledgments

The work is partially supported by the National Natural Science Foundation of China (No. 60970061, No. 61075056, and No. 61103067), the Opening Project of Shanghai Key Laboratory of Digital Media Processing and Transmission (No.2011KF03), and the Fundamental Research Funds for the Central Universities.