

A Novel Progressive Multi-label Classifier for Class-incremental Data

Mihika Dave*, Sahil Tapiawala

Dept. of Electrical & Electronics Engineering
Birla Institute of Technology & Science (BITS)-Pilani
Goa, India

Meng Joo Er, Rajasekar Venkatesan

School of Electrical & Electronics Engineering
Nanyang Technological University
Singapore

Abstract— In this paper, a progressive learning algorithm for multi-label classification to learn new labels while retaining the knowledge of previous labels is designed. New output neurons corresponding to new labels are added and the neural network connections and parameters are automatically restructured as if the label has been introduced from the beginning. This work is the first of the kind in multi-label classifier for class-incremental learning. It is useful for real-world applications in applied fields such as robotics where streaming data are available and the number of labels is often unknown. Based on the Extreme Learning Machine framework, a novel universal classifier with plug and play capabilities for progressive multi-label classification is developed. Experimental results on various benchmark synthetic and real datasets validate the efficiency and effectiveness of our proposed algorithm.

Keywords—extreme learning machines; multi-label classification; on-line learning

I. INTRODUCTION

In the modern context of machine intelligence, the growing importance of classification has motivated the development of several algorithms which are scalable [1][2]. A number of tasks right from the identification of objects in images to the study of emotion-associated brainwaves belong to this category. Classification, which is the identification of the target categories a data sample belongs to, can be divided into single-label and multi-label classification. Single-label classification involves binary and multi-class classification where a data sample is associated with one label only. On the contrary, multi-label classification involves data samples which are simultaneously associated with multiple labels. The learning algorithms are generally of two types: Batch learning and Sequential learning. Batch learning requires pre-collection of training data, and the parameters of the network are calculated by processing all the training data concurrently. While, in online/sequential learning algorithms, the network parameters are updated as and when a new training data arrives [3], [4].

Multi-label classification has become significant due to its rapidly increasing application areas such as text categorization [6]-[8], bioinformatics [9], [10], medical diagnosis [11], scene classification [12], genomics, map labeling [13], marketing, multimedia, music categorization, etc. The rising significance of multi-label classification has spurred a recent growth in its theoretical analysis [14], [15] and development of algorithms for practical applications [16], [17].

The existing multi-label classifiers once trained to classify a specific number of labels, cannot learn new labels without retraining all the labels anew again. Such classifiers work well with the pre-known dataset, but they may not be appropriate for applications such as cognitive robotics or real-time applications of big data where the nature of training data is unknown. For such real-time data, the learning technique must be self-adapting to suit the dynamic needs. Class-incremental Extreme Learning Machines (ELM) has been proposed for multi-class classification [18] but there is no significant work done for multi-label classification. To overcome this shortcoming, a novel learning paradigm based on Extreme Learning Machine is proposed, called the “*progressive-ELM multi-label classifier*” (Pro-EMLC). This is the first method for incremental learning in multi-label classification. It has been successfully tested on benchmark datasets like Scene, Corel5k, and Medical. The promising results we obtained validate the efficacy of our approach. Next section gives a brief description of ELM and Online Sequential-ELM. Section 3 describes the proposed method. Section 4 presents the results obtained and section 5 states the conclusion.

II. BRIEF OVERVIEW OF ELM AND OS-ELM

ELM considers a ‘generalized’ Single Hidden Layer Feedforward Neural Network (SLFN) architecture consisting of n input neurons and m output neurons, with N hidden layer neurons in the second layer. Most neural networks are considered to be universal classifiers or function approximators [19], [20] when all the parameters of the neural network are tuned. It has been previously shown that ELM works for the SLFN architecture without tuning the hidden layer parameters (feature mapping parameters) [21]. We discuss the ELM paradigm in brief in the following paragraphs.

Consider there are N hidden layer neurons, n is the number of features and m is the number of labels. The training data of size N samples is of the form $\{(\mathbf{x}_i, \mathbf{y}_i) | \mathbf{x}_i \in R^n, \mathbf{y}_i \in R^m, i = 1, \dots, N\}$, where \mathbf{x}_i represents input feature vector and \mathbf{y}_i represents the output label vector. Label space $L = \{Y_1, Y_2, \dots, Y_M\}$. The predicted output of SLFN ‘ o_j ’ for $j = 1, \dots, N$ is:

$$\sum_{i=1}^N \beta_i g(\mathbf{w}_i \cdot \mathbf{x}_j + b_i) = o_j \quad (1)$$

Where, $g(x)$ is the activation function, $\mathbf{w}_i = [w_{i1}, w_{i2}, \dots, w_{in}]^T$ is the input weight, $\beta_i = [\beta_{i1}, \beta_{i2}, \dots, \beta_{im}]^T$ is the output weight, and b_i is the hidden layer bias.

In ELM, the input weights and hidden layer bias are assigned randomly. To minimize the difference between the actual and predicted output, β_i should be found such that output class is equal to the target class.

$$\sum_{j=1}^{N'} \|\mathbf{o}_j - \mathbf{y}_j\| = 0 \quad (2)$$

Thus, ELM output network can be written as

$$\sum_{i=1}^{N'} \beta_i g(\mathbf{w}_i \cdot \mathbf{x}_j + b_i) = \mathbf{y}_j, \quad j=1, \dots, N \quad (3)$$

In matrix form,

$$\mathbf{H}\boldsymbol{\beta} = \mathbf{Y} \quad (4)$$

Where,

$$\mathbf{H} = \begin{bmatrix} g(\mathbf{w}_1 \cdot \mathbf{x}_1 + b_1) & \dots & g(\mathbf{w}_{N'} \cdot \mathbf{x}_1 + b_{N'}) \\ \vdots & \ddots & \vdots \\ g(\mathbf{w}_1 \cdot \mathbf{x}_N + b_1) & \dots & g(\mathbf{w}_{N'} \cdot \mathbf{x}_N + b_{N'}) \end{bmatrix}_{N \times N'}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_1^T \\ \vdots \\ \beta_{N'}^T \end{bmatrix}_{N' \times m} \quad \text{and} \quad \mathbf{Y} = \begin{bmatrix} y_1^T \\ \vdots \\ y_N^T \end{bmatrix}_{N \times m}$$

Using the Moore-Penrose generalized inverse (\mathbf{H}^+) of hidden layer matrix \mathbf{H} , we can get the output feature mapping matrix $\boldsymbol{\beta}$ of the ELM network.

$$\boldsymbol{\beta} = \mathbf{H}^+ \mathbf{Y} \quad (5)$$

Where,

$$\mathbf{H}^+ = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$$

In real time applications, the complete training data is seldom available and as the data arrives sequentially, the ELM has to be retrained with a combination of the new data and the previously available data. There is no provision to retain the learning from the previous data and train the network for only the new data. On-line Sequential Extreme Learning Machine (OS-ELM) retains the knowledge of the previous training data and can learn data over previously available data chunk-by-chunk with varying chunk sizes [22]. The above-mentioned algorithm of ELM is extended to OS-ELM as described below.

Let the initial block of training data have N_0 samples. For the initial block, compute

$$\mathbf{M}_0 = (\mathbf{H}_0^T \mathbf{H}_0)^{-1} \quad (6)$$

$$\boldsymbol{\beta}_0 = \mathbf{M}_0 \mathbf{H}_0^T \mathbf{Y}_0 \quad (7)$$

Where $\mathbf{H}_0 = [\mathbf{h}_1 \dots \mathbf{h}_{N_0}]^T$ is the feature mapping (hidden layer) matrix computed as above and $\mathbf{Y}_0 = [y_1, \dots, y_{N_0}]^T$ output vector for N_0 samples

For the incoming $(k+1)^{th}$ sequential data, Recursive Least Squares (RLS) approximation is used to retain the learning. Updating can be done as,

$$\mathbf{M}_{k+1} = \mathbf{M}_k - \frac{\mathbf{M}_k \mathbf{h}_{k+1} \mathbf{h}_{k+1}^T \mathbf{M}_k}{1 + \mathbf{h}_{k+1}^T \mathbf{M}_k \mathbf{h}_{k+1}} \quad (8)$$

$$\boldsymbol{\beta}_{k+1} = \boldsymbol{\beta}_k + \mathbf{M}_{k+1} \mathbf{h}_{k+1} (\mathbf{Y}_{k+1}^T - \mathbf{h}_{k+1}^T \boldsymbol{\beta}_k) \quad (9)$$

Where $k = 0, 1, 2, \dots, N-N_0-1$.

The calculated value of $\boldsymbol{\beta}$ is used for predicting the output matrix. The detailed theory and mathematics involved in ELM and OS-ELM are given in [22] [23].

III. PROPOSED METHOD

The steps of the proposed 'progressive-ELM multi-label classifier' are:

1) *Processing of input*: The raw input data is processed so that the output label, corresponding to each input sample, is an m-tuple with -1 or 1 representing the belongingness to each of the labels in the label space L.

2) *Initialization*: The input weights and the hidden layer bias are assigned at random in accordance with the ELM learning paradigm. The number of hidden layer neurons N' is fixed such that over-fitting does not occur. The optimal value N' is found by carrying out different epochs with varying number of hidden neurons, plotting it w.r.t. the training and cross-validation accuracy and choosing the optimal number of hidden neurons.

3) *ELM training – Initial learning*: The hidden layer output matrix \mathbf{H}_0 is calculated for an initial block of N_0 data samples.

$$\mathbf{H}_0 = [\mathbf{h}_1, \dots, \mathbf{h}_{N_0}]^T \quad (10)$$

Where $\mathbf{h}_i = [g(\mathbf{w}_1 \cdot \mathbf{x}_i + b_1), \dots, g(\mathbf{w}_{N'} \cdot \mathbf{x}_i + b_{N'})]^T, i = 1, 2, \dots, N_0$.

Using \mathbf{H}_0 , the initial values of \mathbf{M}_0 and $\boldsymbol{\beta}_0$ are estimated as

$$\mathbf{M}_0 = (\mathbf{H}_0^T \mathbf{H}_0)^{-1} \quad (11)$$

$$\boldsymbol{\beta}_0 = \mathbf{M}_0 \mathbf{H}_0^T \mathbf{Y}_0 \quad (12)$$

4) *ELM training – Sequential Learning*: The subsequent data arriving at the network can be trained one-by-one or chunk-by-chunk. Let the chunk size be 'b'. For b=1, data is trained one-by-one. The incoming data may have the presence of a new label(s), indicated by the increase in tuple size. Presence/Absence of a new label is handled as below:

a) *Absence*: In the absence of a new label, no special computations are required and step 5 is executed directly.

b) *Presence*: In the presence of a new label, the network is to be recalibrated to accommodate a new label, while retaining the old knowledge. Let 'c' new labels be introduced and m labels of data are currently learnt. Introducing 'c' new labels at any instant k+1 modifies the dimensions of the output weight matrix $\boldsymbol{\beta}_{N' \times m}$ to $\boldsymbol{\beta}_{N' \times m+c}$. Transformed output weight matrix $\bar{\boldsymbol{\beta}}_k$ is given as $\bar{\boldsymbol{\beta}}_k = (\boldsymbol{\beta}_k)_{N' \times m} \mathbf{I}_{m \times m+c}$ where, $\mathbf{I}_{m \times m+c}$ is m X m+c dimensional rectangular identity matrix.

$$\bar{\boldsymbol{\beta}}_{k, N' \times m+c} = (\boldsymbol{\beta}_k)_{N' \times m} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \end{bmatrix}_{m \times m+c} \quad (13)$$

$$\bar{\boldsymbol{\beta}}_{k, N' \times m+c} = [(\boldsymbol{\beta}_k)_{N' \times m} \quad \mathbf{O}_{N' \times c}]_{N' \times m+c} \quad (14)$$

where $\mathbf{O}_{N' \times c}$ is a zero matrix.

The output values corresponding to new labels is -1 for previously learnt samples. Therefore, the k-learning step update for the 'c' new labels ($(\Delta \boldsymbol{\beta}_k)_{N' \times c} (\Delta \boldsymbol{\beta}_k)_{N' \times c}$) can be expressed as,

$$(\Delta \boldsymbol{\beta}_k)_{N' \times c} = (\mathbf{M}_k)_{N' \times N'} (\mathbf{h}_k^T)_{N' \times b} \begin{bmatrix} -1 & \dots & -1 \\ \vdots & \ddots & \vdots \\ -1 & \dots & -1 \end{bmatrix}_{b \times c} \quad (15)$$

$$(\Delta\beta_k)_{N'Xc} = -(\mathbf{M}_k)_{N'XN'} (\mathbf{h}_k^T)_{N'Xb} J_{bXc} \quad (16)$$

Where, J_{bXc} is an all-ones matrix of size $b \times c$.

$$(\Delta\tilde{\beta}_k)_{N'Xm+c} = [O_{N'Xm} \quad -(\mathbf{M}_k)_{N'XN'} (\mathbf{h}_k^T)_{N'Xb} J_{bXc}] \quad (17)$$

The recalibrated output weight matrix $(\beta_k)_{N'X(m+c)}$ is calculated as,

$$(\beta_k)_{N'X(m+c)} = \tilde{\beta}_{kN'Xm+c} + (\Delta\tilde{\beta}_k)_{N'Xm+c} \quad (18)$$

Upon simplification, $(\beta_k)_{N'X(N'+c)}$ can be expressed as,

$$(\beta_k)_{N'X(m+c)} = [(\beta_k)_{N'Xm} \quad (\Delta\beta_k)_{N'Xc}] \quad (19)$$

The hidden layer output vector h_{k+1} is calculated. The output weight is updated according to the Recursive Least Square algorithm [22],

$$\mathbf{M}_{k+1} = \mathbf{M}_k - \mathbf{M}_k \mathbf{h}_{k+1}^T (\mathbf{I} + \mathbf{h}_{k+1} \mathbf{M}_k \mathbf{h}_{k+1}^T)^{-1} \mathbf{h}_{k+1} \mathbf{M}_k \quad (20)$$

$$\beta_{k+1} = \beta_k + \mathbf{M}_{k+1} \mathbf{h}_{k+1} (\mathbf{Y}_{k+1}^T - \mathbf{h}_{k+1}^T \beta_k) \quad (21)$$

5) *ELM testing and thresholding*: Raw output matrix Y of the testing data samples is calculated using $Y = H\beta$. Since the number of labels a samples belongs to is unknown, thresholding of raw output is suggested. Thresholding refers to the application of a threshold value, based on the separation between two categories of labels (Labels that the data sample belongs to and labels the data sample does not belong to), to identify the number of labels and the target class labels corresponding to the input data sample. As a trivial case, zero threshold is chosen. Passing the raw output through bipolar step function gives the samples' association to the label(s).

Thus, the proposed algorithm allows for class-incremental training of multi-label data.

IV. EXPERIMENTAL RESULTS

The proposed method is tested on benchmark datasets, namely Scene, Core15k and Medical, having a varied number of labels (6 to 374) and belonging to diverse domains like multimedia and text. The degree of 'multi-labelness' of datasets is measured by label density (LD) and label cardinality (LC) as defined by Tsoumakas et al [24]. Label density of the tested datasets ranges from 0.009 to 0.178 and label cardinality ranges from 1.07 to 3.52. The characteristics of these datasets are given

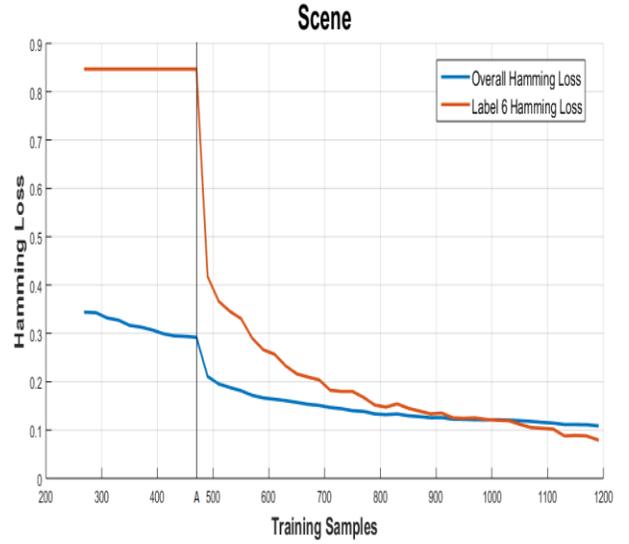
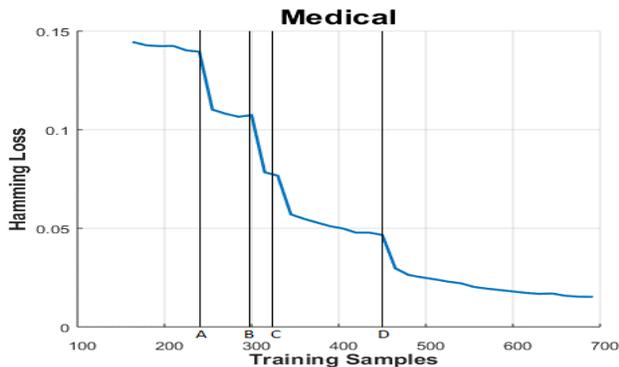


Fig. 1. Hamming Loss for Scene Dataset

TABLE I. CHARACTERISTICS OF DATASET

Characteristic	Dataset		
	Scene	Core15k	Medical
Domain	Multimedia	Multimedia	Text
No. of features	294	499	1449
No. of samples	2407	5000	978
No. of labels	6	374	45
Label Density	0.178	0.009	0.027
Label Cardinality	1.07	3.52	1.25

in Table I. To perform testing for progressive learning, the data samples are redistributed such that the number of labels present in the initial block of data is less than the total number of labels present in the dataset. New labels can be introduced in the streaming data one-by-one or in groups.

The learning curve or the hamming loss is plotted for Scene dataset (having 6 labels) in Fig. 1. The hamming loss decreases with increasing number of samples. Till point A, when only five

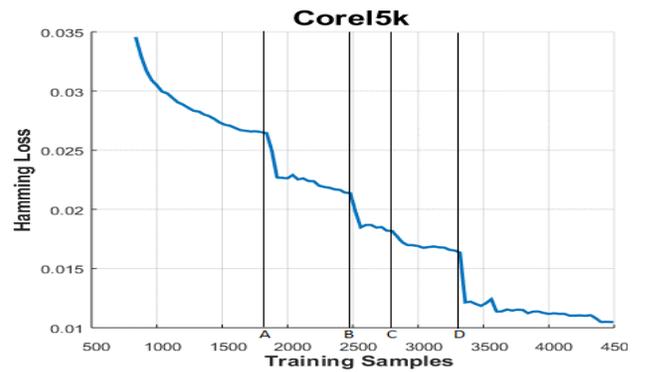


Fig. 2. Hamming loss for medical (39+2+2+1+1) and core15k (368+2+1+1+2) dataset

TABLE II. PERFORMANCE METRIC

Data set	LIP	H	Acc	Pre	Rec	F1	T1	T2
Scene	5+1	0.104	0.609	0.627	0.659	0.643	2.266	0.042
Scene	4+1+1	0.139	0.569	0.584	0.699	0.636	2.291	0.080
Medical	44+1	0.012	0.693	0.740	0.729	0.734	0.604	0.025
Medical	39+2+2+1+1	0.016	0.653	0.695	0.731	0.712	0.611	0.034
Medical	38+3+4	0.023	0.585	0.622	0.739	0.675	0.614	0.027
Corel5k	373+1	0.010	0.057	0.164	0.059	0.087	5.348	0.044
Corel5k	368+2+1+1+2	0.010	0.055	0.151	0.062	0.088	5.396	0.056

(LIP- Label Introduction Pattern, H- Hamming Loss, Acc- Accuracy, Pre- Precision, Rec- Recall, F1- F1 Score, T1- Train Time, T2- Test Time)

labels are learnt, the overall hamming loss is higher since the prediction for the 6th label is incorrect. After introducing the remaining label (label 6), at point A, in the sequential phase, there is a sharp decrease in hamming loss suggesting that the prediction for label 6 has improved. The hamming loss for the label 6 alone is also shown in Fig. 1. The learning curves for Medical and Corel5k datasets are shown in Fig. 2. For the Medical dataset, 39 out of 45 labels were introduced initially. To verify the behavior towards incremental labels, two labels are introduced at point A, two labels at point B, one label at point C, and finally the last label at point D for the Medical dataset. The falling hamming loss is evidence of the improvement in the

TABLE III. STATE-OF-THE-ART MULTI-LABEL CLASSIFIERS

Method Name	Category
Classifier Chain (CC)	SVM
QWeighted approach for Multi-label Learning (QWML)	SVM
Hierarchy Of Multi-label ClassifierS (HOMER)	SVM
Multi-Label C4.5 (ML-C4.5)	Decision Trees
Predictive Clustering Trees (PCT)	Decision Trees
Multi-Label k-Nearest Neighbors (ML-kNN)	Nearest Neighbors
Ensemble of Classifier Chains (ECC)	SVM
Random Forest Predictive Clustering Trees (RF-PCT)	Decision Trees
Random Forest of ML-C4.5 (RFML-C4.5)	Decision Trees

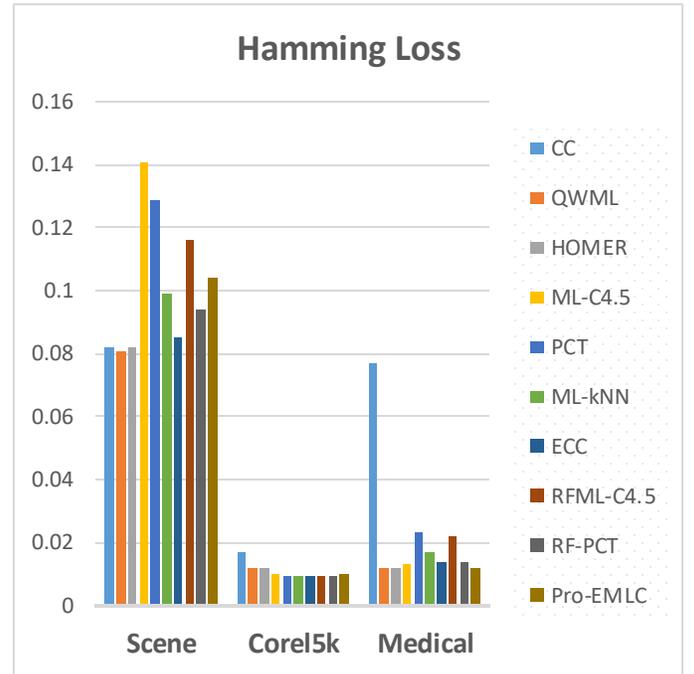


Fig. 3. Comparison of hamming loss with other multi-label classifiers

learning resulting in better prediction. This pattern of introducing new labels is represented as 39+2+2+1+1 in Table II. Similar protocol was followed for Corel5k dataset with the introduction pattern being 368+2+1+1+2.

A 10-fold cross-validation is carried out for various combinations of label introduction in Scene (6 labels), Medical (45 labels) and Corel5k (374 labels) data sets. Various measures to evaluate the performance of multi-label classifiers, as explained by Tsoukamas et al [24] are calculated for the proposed algorithm. The results are presented in Table II. Label introduction pattern is the number of labels introduced in the initial phase followed by the number of those in the sequential phase.

The proposed algorithm was compared to the state-of-the-art online learning multi-label algorithms mentioned in Table III.

TABLE IV. COMPARISON OF TRAINING TIME FOR VARIOUS MULTI-LABEL CLASSIFIERS

Classifier	Dataset		
	Scene	Medical	Core15k
CC	99	1125	28
QWML	195	2388	40
HOMER	68	771	16
ML-C4.5	8	369	3
PCT	2	30	0.6
ML-kNN	14	389	1
ECC	319	10073	103
RFML-C4.5	10	385	7
RF-PCT	23	902	27
Pro-EMLC	2.266	5.396	0.611

Despite being an online-class incremental algorithm, the hamming loss of the proposed Progressive - ELM Multi-label Classifier (Pro-EMLC), as shown in Fig. 3, was comparable to the state-of-the-art multi-label classification techniques for batch learning described by Madjarov et al. [25], for all the tested datasets. A similar comparison of training time was also made as given in Table IV. Since the algorithm is based on ELM, the training and testing speed is very high, making it highly suitable for real-time applications of big data analysis. This suggests that the performance of the algorithm is superior for a class incremental algorithm and hence can be used widely.

V. CONCLUSIONS

The proposed progressive-ELM multi-label classifier is the first of its kind. It can be used for online/streaming big data applications with known number of labels as well as for real-time applications such as cognitive robotics where the number of labels is unknown. It has shown high speed and high performance metric for all the 3 tested benchmark datasets. Based on these promising results, the proposed ELM-based classifier can be considered fit for achieving high performance with speed for multi-label classification, especially in incremental learning areas.

ACKNOWLEDGEMENT

The authors would like to acknowledge the funding support from the Ministry of Education, Singapore (Tier 1 AcRF, RG30/14).

REFERENCES

[1] G. B. Huang, Z. Hongming, D. Xiaojian and Z. Rui, "Extreme Learning Machine for Regression and Multiclass Classification," *IEEE Transactions on Systems, Man and Cybernetics, Part B*, vol. 42, no. 2, pp. 513-529, 2012

[2] K. Crammer and C. Gentile, "Multiclass classification with bandit feedback using adaptive regularization," *Machine Learning*, vol. 90, no. 3, pp. 347-383, 2013.

[3] M. Pratama, S. G. Anavatti, and J. Lu. Recurrent classifier based on an incremental metacognitive-based scaffolding algorithm. *IEEE Transactions on Fuzzy Systems*, 23(6):2048-2066, 2015.

[4] M. Pratama, J. Lu, and G. Zhang. Evolving Type-2 Fuzzy Classifier. *IEEE Transactions on Fuzzy Systems*, PP(99):1-1, 2015.

[5] G. Tsoumakas, I. Katakis and I. Vlahavas. Mining Multi-label Data. In O. Maimon and L. Rokach, editors, *Data Mining and Knowledge Discovery Handbook*, pages 667-685, Springer US, 2010.

[6] X. Luo and A. N. Zincir-Heywood. Evaluation of Two Systems on Multi-class Multi-label document classification. In Hacid M.S., Murray N.V., Ras Z.W. and Tsumoto S, editors, *Foundations of Intelligent Systems*, volume 3488 of *Lecture Notes in Computer Science*, pages 161-169. Springer Berlin Heidelberg, 2005.

[7] D. Tikk and G Biro. Experiments with multi-label text classifier on the Reuters collection. In *Proceedings of the international conference on computational cybernetics (ICCC 03)*, pages 33-38, 2003.

[8] K. Yu, S. Yu and V. Tresp. Multi-label informed latent semantic indexing In *Proceedings of the 28th annual international ACM SIGIR conference on Research and Development in information retrieval, SIGIR'05*, pages 258-265, New York, NY, USA, 2005.

[9] A. Elisseeff and J. Weston. A kernel method for multi-labelled classification. In *Advances in Neural Information Processing Systems*, pages 681-687, 2001.

[10] M. L. Zhang and Z. H. Zhou, A k-nearest neighbour based algorithm for multi-label classification. In *2005 IEEE International Conference on Granular Computing*, volume 2, pages 718-721, 2005.

[11] A. Karalic and V. Pirnat. Significance level based multiple tree classification. *Informatica*, 15(5):12, 1991.

[12] M. Boutell, X. Shen, J. Luo and C. Brouwn. Multi-label semantic scene classification. *Technical Report, Department of Computer Science, University of Rochester, USA*, 2003.

[13] B. Zhu and C. K. Poon. Efficient Approximation Algorithms for Multi-label Map Labelling. In *Algorithms and Computation*, volume 1741 of *Lecture Notes in Computer Science*, pages 143-152. Springer Berlin Heidelberg, 1999.

[14] K. Dembczy, W. Waegeman, W. Cheng and E. Hullemeier, "On Label dependence and Loss Minimization in Multilabel Classification," *Machine Learning*, vol. 88, no. 2, pp. 5-45, 2012

[15] W. Gao and Z.-H. Zhou, "On the consistency of Multi-label Classification," *Artificial Intelligence*, vol. 199, pp. 22-44, 2013

[16] J. Petterson and T. Caetano, "Submodular Multi-label Learning," in *Advances in Neural Information Processing Systems*, Granada, Spain, 2011.

[17] H.-F. Yu, P. Jain, P. Kar and I. Dhillon, "Large-scale multi-label learning with missing labels," in *Proceedings of the 31st International Conference on Machine Learning*, Beijing, China, 2014.

[18] Z. Zhao, Z. Chen, Y. Chen, S. Wang, and H. Wang. A class incremental extreme learning machine for activity recognition. *Cognitive Computation*, 6(3):423-431, 2014.

[19] K. Hornik, "Approximation capabilities of multilayer feedforward networks," *Neural Network*, vol. 4, pp. 251-257, 1991.

[20] M. Leshno, V. Y. Lin, A. Pinkus and S. Schocken, "Multilayer feedforward networks with a nonpolynomial activation function can approximate any function," *Neural Network*, vol. 6, no. 6, pp. 861-867, 1993.

[21] G.-B. Huang, L. Chen and C.-K. Siew, "Universal approximation using incremental constructive feedforward networks with random hidden nodes," *IEEE Transactions on Neural Networks*, vol. 17, no. 4, pp. 879-892, 2006.

[22] N. Y. Liang, G. B. Huang, P. Saratchandran and N. Sundararajan. A Fast and Accurate Online Sequential Learning Algorithm for Feedforward Networks. *IEEE Transactions on Neural Networks*, 17(6):1411-1423, 2006.

[23] G. B. Huang, D. H. Wang and Y. Lan. Extreme Learning Machines: A Survey. *International Journal of Machine Learning and Cybernetics*, 2(2):107-122, 2011.

[24] G. Tsoumakas and I. Katakis. Multi-label Classification: An Overview, *International Journal of Data Warehousing and Mining*, 3(3): 1-13, 2007.

[25] G. Madjarov, D. Kocev, D. Gjorgjevikj and S. Dzeroski. An extensive experimental comparison of methods for multi-label learning. *Pattern Recognition*, 45(9): 3084-3104, 2012.