

LABEL PROPAGATION ADAPTIVE RESONANCE THEORY FOR SEMI-SUPERVISED CONTINUOUS LEARNING

Taehyeong Kim^{1,2} Injune Hwang¹ Gi-Cheon Kang¹ Won-Seok Choi¹
Hyunseo Kim¹ Byoung-Tak Zhang¹

¹Seoul National University, Seoul, Republic of Korea

²AI Lab, CTO Division, LG Electronics, Seoul, Republic of Korea

{thkim, ijhwang, gckang, wchoi, hskim, btzhang}@bi.snu.ac.kr

ABSTRACT

Semi-supervised learning and continuous learning are fundamental paradigms for human-level intelligence. To deal with real-world problems where labels are rarely given and the opportunity to access the same data is limited, it is necessary to apply these two paradigms in a joined fashion. In this paper, we propose *Label Propagation Adaptive Resonance Theory* (LPART) for semi-supervised continuous learning. LPART uses an online label propagation mechanism to perform classification and gradually improves its accuracy as the observed data accumulates. We evaluated the proposed model on visual (MNIST, SVHN, CIFAR-10) and audio (NSynth) datasets by adjusting the ratio of the labeled and unlabeled data. The accuracies are much higher when both labeled and unlabeled data are used, demonstrating the significant advantage of LPART in environments where the data labels are scarce.

Index Terms— Label propagation, adaptive resonance theory, semi-supervised learning, continuous learning

1. INTRODUCTION

Over the last few years, the deep neural network models have shown remarkable progresses especially in visual object recognition, speech recognition and autonomous robot control. However, using deep learning has a major practical shortcoming that it requires time and labor not only to collect a massive amount of data but also to label them. In this aspect, the fields of semi-supervised learning, continuous learning, transfer learning and meta-learning have been in the spotlight.

The semi-supervised learning paradigm [1] tackles the problem in which the unlabeled data is abundant but the labeled data is extremely limited. The continuous learning paradigm [2] aims to learn without catastrophic forgetting of the former knowledge from sequential data, allowing the model to adapt to the ever-changing environments. Although

each paradigm is promising on its own, we argue that both paradigms should be applied in a joined fashion to deal with many real-world problems, where labeled data is rarely given and the data once learned is hard to access again.

In this study, we propose *Label Propagation Adaptive Resonance Theory* (LPART) for semi-supervised continuous learning. Adaptive Resonance Theory (ART) is a solution to continuous learning inspired by brain information processing mechanisms [3], and various label propagation methods have been studied for semi-supervised learning [4, 5]. However, these label propagation methods are not suitable for the environment with limited opportunity to access the same data because they require the data repeatedly to be learned.

Therefore, we propose an online label propagation method for continuous learning in the ART network. Specifically, we use a two-fold learning process: (1) feature extraction using variational autoencoders (VAE) [6] and (2) clustering of the extracted features and inference of the classes using LPART. First, we train VAE in a weakly-supervised manner by using the *pair loss*. Then, LPART takes the features as input and learns to infer the classes of unlabeled nodes by leveraging the label propagation method. If the amount of labeled data is not enough, the inference of the unlabeled node could be inaccurate. Therefore, we use the metrics to measure the uncertainties so that LPART could defer its classification decisions for the nodes with high uncertainty. We experimentally confirmed that our proposed model is able to learn continuously without catastrophic forgetting, even with the rarely labeled data.

2. ADAPTIVE RESONANCE THEORY

ART is a self-organizing neural network inspired by the brain information processing mechanisms. ART uses the interaction of ‘top-down’ expectation and ‘bottom-up’ sensory information to learn adaptively, using resonance. There are two general principles in ART: (1) The knowledge is strengthened when the sensation is strong enough and the expectation matches well with the sensation, and (2) if there is no expect-

This work was partly supported by the Korea government (2019-0-01367-BabyMind, 2015-0-00310-SW.StarLab, 2017-0-01772-VTT, 2018-0-00622-RMI, P0006720-GENKO).

tation that matches the sensation, new knowledge is learned. In terms of being conservative while learning new, the ART system can be a solution for the continuous learning.

Various ART networks such as Fuzzy ART [7], ARTMAP [8], and Fuzzy ARTMAP (FAM) [9] have been studied. Fuzzy ART can process real-valued data by using fuzzy set theory. ARTMAP uses supervised learning and classification system that is built up from a pair of ART modules. FAM integrates the advantages of both Fuzzy ART and ARTMAP.

There are also other ART networks for learning in a semi-supervised manner, such as semi-supervised Bayesian ARTMAP (SSBA) [10] and semi-supervised Fuzzy ARTMAP (ssFAM) [11]. SSBA employs EM algorithm based on Bayesian ARTMAP (BA) [12] to adjust its parameters, which realizes the soft assignment of training samples instead of the winner-take-all strategy. ssFAM relies on FAM but adopts a tunable network parameter called *category prediction error tolerance*, which achieves semi-supervised learning.

The present study differs from previous studies in that it applies an online label propagation mechanism for semi-supervised continuous learning. In addition, the label propagation mechanism can be applied to various kinds of ART networks, which allows the extension from LPART. Also, the uncertainty measurement methods proposed in this study can be used to filter reliable classification results.

3. METHODS

3.1. Feature Extraction with VAE

It is necessary to extract features that are easy to cluster for ART to classify high-dimensional data properly. VAE, a deep learning-based unsupervised learning method, is widely used to extract useful features from data [13, 14]. However, basic VAE does not have any explicit constraints to improve clustering. In this regard, the study of learning representations using *oracle triplets* provides the insights needed for this study [15]. We also repurposed the VAE architecture for semi-supervised continuous learning.

In this study, we use a simplified triplet-based VAE to extract features. It uses only some dimensions d in the latent space for VAE encoder to produce class-embedded representation μ_d . Additional *pair loss* is introduced, which depends on whether or not the class of the previous sample and the current sample are the same. The *pair loss* between previous and current class-embedded representations is defined using the L2 distance as a similarity measure (Equation 1).

$$\mathcal{L}_{pair} = \begin{cases} \|\mu_{d,prev} - \mu_{d,curr}\|_2, & y_{prev} = y_{curr} \\ -\|\mu_{d,prev} - \mu_{d,curr}\|_2, & otherwise \end{cases} \quad (1)$$

Here, y denotes a label of the input sample. We optimize parameters by maximizing the ELBO (evidence lower bound) [6] and minimizing the *pair loss*. With a scaling factor λ , the

total loss to be minimized is as shown in Equation 2.

$$\mathcal{L} = -ELBO + \lambda\mathcal{L}_{pair} \quad (2)$$

3.2. The LPART Algorithm

When an input data x_i is given, we encode it using the VAE previously described to get the 0-to-1 normalized class-embedded representation as r_i . As in Fuzzy ART, we also compute the complement coding I_i of r_i . For a node j with a weight vector w_j , the choice function T_j and the match function V_j of I_i are defined as:

$$T_j(I_i) = \frac{\|I_i \wedge w_j\|_1}{\alpha + \|w_j\|_1}, \quad V_j(I_i) = \frac{\|I_i \wedge w_j\|_1}{\|I_i\|_1} \quad (3)$$

where \wedge is the element-wise minimum operator, $\alpha > 0$ is the choice parameter and $\|\cdot\|_1$ denotes the L1 norm of a vector.

If the value of $V_j(I_i)$ is greater than a vigilance parameter ρ , we say that the node j has matched x_i or been *activated*. Among all of the activated nodes, a *winner* J with the highest value of T_j is selected. It can be seen as the best-fit node for the input, and we update its weight vector considering I_i with a learning rate β between 0 and 1 as shown in Equation 4.

$$w_j^{new} = \beta(I_i \wedge w_j^{old}) + (1 - \beta)w_j^{old} \quad (4)$$

On the other hand, if no node matches the input, a new node is created with an initial parameter set as I_i . This newborn node can grow larger throughout the subsequent iterations. The creation of a node is more frequent with a larger value of vigilance parameter ρ . By manipulating the value of ρ , we can balance the rigidity of the node. If it gets too small, one node covers up too many inputs, making the consistency of the node vague. Therefore, we use sufficiently large value for the vigilance parameter.

Another crucial part of LPART, the label propagation mechanism, will be explained in the following section. The overall LPART algorithm is described in Algorithm 1.

3.3. Label Propagation Mechanism

Label propagation, a mechanism for semi-supervised learning, is a method of inferring a class of unlabeled cluster with the help of labeled ones [16]. It assumes that the clusters close to each other in the feature space tend to belong to similar classes. In LPART, label propagation is triggered when an input data activates two or more nodes. The co-activated nodes can be considered to be located in the vicinity of each other in the feature space, which in turn implies a high relevance between them. It is natural, therefore, to estimate the label of a label-absent node—a node that does not contain any input with a label—using the labels of co-activated nodes. The numerical value of the label itself is meaningless, so we use a distribution over all labels instead of a single value. We call this a *label density function* and denote by q , where $q_j(y)$

Algorithm 1: The LPART algorithm.

```
1 for  $x_i, y_i$  in dataset do //  $y_i$  can be absent
2    $r_i \leftarrow \text{Encode}(x_i)$ 
3    $I_i \leftarrow [r_i, \vec{1} - r_i]$  // concatenation
4    $A \leftarrow \{\}$ 
5   for  $j$  in  $1, \dots, N$  do //  $N$  is the number of nodes
6      $T_j \leftarrow \|I_i \wedge w_j\|_1 / (\alpha + \|w_j\|_1)$ 
7      $V_j \leftarrow \|I_i \wedge w_j\|_1 / \|I_i\|_1$ 
8     if  $V_j \geq \rho$  then
9       if  $y_i$  is given then
10         $q_j(y_i) \leftarrow q_j(y_i) + 1$ 
11         $A \leftarrow A \cup \{j\}$ 
12   if  $A$  is not empty then
13     LabelPropagate( $A$ ) // if  $|A| > 1$ 
14      $J \leftarrow \arg \max_{j \in A} (T_j)$ 
15      $w_J \leftarrow \beta(I_i \wedge w_J) + (1 - \beta)w_J$ 
16   else
17     CreateNode( $x_i, y_i$ ) //  $q_n(y_i) \leftarrow 1$ 
```

roughly means how probable a node j will be in class y . When a new node n is created, q_n is initialized to the zero vector.

Once a labeled sample is added to a node, the density of its label increases by one. For a label-absent node k , label density function is updated by averaging those of co-activated nodes:

$$q_k^{new}(y) = \left(\delta \times \frac{\sum_{j \in A - \{k\}} q_j^{old}(y)}{\sum_{y'} \sum_{j \in A - \{k\}} q_j^{old}(y')} + (1 - \delta) \times \frac{q_k^{old}(y)}{\sum_{y'} q_k^{old}(y')} \right) \times \frac{1}{C} \quad (5)$$

where δ is a propagation rate. The reason why the sum of q_k over all labels is less than one is to indicate that it is still not certain which class this node belongs to. $C > 1$ can be interpreted as a kind of uncertainty parameter, which will be further discussed in Section 3.4.

Finally, the probability distribution of labels for each node is easily obtained by normalizing the label density function:

$$p_j(y) = \frac{q_j(y)}{\sum_{y'} q_j(y')} \quad (6)$$

and we can infer a class of input data by finding a winner node and then selecting a label with the highest probability.

3.4. Measurement of Uncertainty

We use two different metrics to measure the uncertainty of the classification results. The first uncertainty, $u^1(x_i)$, is measured by the entropy of the classification probability [17], as shown in Equation 7. This method measures how evenly distributed the categories of labeled data learned by each node,

which can be used to identify high-impurity nodes.

$$u^1(x_i) = - \sum_y p_{J(x_i)}(y) \log p_{J(x_i)}(y) \quad (7)$$

$J(x_i)$ is a winning node with given input x_i and $p_{J(x_i)}(y)$ is the probability that this node belongs to class y .

The second uncertainty $u^2(x_i)$ is based on the number of labeled input data in each node as shown in Equation 8. It allows them to filter out highly unreliable recognition results from the datasets with insufficient number of labels.

$$u^2(x_i) = 1 - \tanh(k \cdot \sum_y q_{J(x_i)}(y)) \quad (8)$$

Here, k is a constant for sensitivity. The combination of these two uncertainties comes in handy for real-world problems; for example, we can withhold judgments on samples with high uncertainties during continuous learning.

4. EXPERIMENTS

4.1. Dataset

We investigated our proposed model using MNIST [18], SVHN [19], CIFAR-10 [20], and NSynth [21] datasets. The MNIST and SVHN datasets consist of the digit images from 0 to 9. The CIFAR-10 dataset contains 60,000 color images in 10 classes such as airplanes, cars, birds and cats. The NSynth dataset contains 305,979 4-second audio recordings from 1,006 instruments. Each recording is labeled with one of the 11 high-level groups such as bass, keyboard, synth lead, and vocal. Note that since we only used 12,678 validation split as training data (with 4,096 original test split), the synth lead included only in the training split was omitted. All audio data was converted to spectrogram images for feature extraction. The features extracted using the VAE are shown in Figure 1.

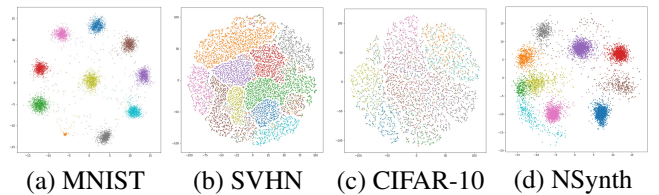


Fig. 1. The visualization of extracted features. For plotting (b) and (c), t-SNE [22] was used to reduce dimensions.

4.2. Semi-supervised Learning

We conducted experiments to evaluate the semi-supervised learning performance of our model on the aforementioned datasets. As shown in Table 1, we verify the semi-supervised classification performance by adjusting the ratio of the labeled and unlabeled data. Also, we compare our proposed model with the FAM. Because the FAM model is based on fully-supervised learning, we could not report the performance of the FAM on semi-supervised settings. We calculated the average performance of 30 trials.

Table 1. Classification accuracy of our model (LPART) compared to FAM trained for a single epoch with various probabilities of the labeled data. The mean and standard deviation are drawn from 30 trials for each experiment. (unit : %)

Dataset rate (labeled, unlabeled)	MNIST ($\rho = 0.99$)		SVHN ($\rho = 0.98$)		CIFAR-10 ($\rho = 0.95$)		NSynth ($\rho = 0.95$)	
	FAM	LPART	FAM	LPART	FAM	LPART	FAM	LPART
0.1%, not used	46.2±3.0	46.5±3.4	48.1±2.8	48.8±2.9	33.0±2.5	30.9±2.9	46.8±9.3	48.3±10.9
0.1%, 99.9%	-	94.2±1.0	-	73.7±1.3	-	39.5±2.0	-	63.1±11.7
0.5%, not used	71.8±1.9	70.9±1.7	60.5±2.0	59.4±1.9	38.8±1.2	37.1±1.5	67.2±4.3	68.3±4.5
0.5%, 99.5%	-	95.0±0.3	-	74.5±0.4	-	42.4±0.9	-	85.2±2.1
1.0%, not used	79.5±1.6	79.8±1.6	63.6±1.7	63.4±1.3	40.3±0.7	38.6±1.3	73.6±3.3	74.0±4.0
1.0%, 99.0%	-	95.3±0.3	-	74.8±0.4	-	42.8±0.7	-	87.6±1.6
5.0%, not used	90.2±0.8	90.1±0.8	70.6±0.8	70.8±0.6	41.9±0.6	42.6±0.9	84.3±1.8	84.4±1.4
5.0%, 95.0%	-	96.2±0.2	-	76.0±0.4	-	44.1±0.5	-	90.0±0.7

4.3. Semi-supervised Continuous Learning

In this experiment, we expanded our experiment to semi-supervised continuous learning. We performed two experiments on NSynth dataset: (1) accuracy comparison between LPART with FAM by epochs, and (2) the uncertain sample rate and classification accuracies using the uncertainty measurement methods by epochs. We also calculated the average performance of 10 trials.

5. RESULTS AND DISCUSSION

5.1. Semi-supervised Learning

The classification accuracies on the four datasets are summarized in Table 1, with different few-labeled data probability settings. In all experimental setups, the best results were obtained from our model using both labeled and unlabeled data, which is much higher than the accuracy using the labeled data only. The performance drops as the amount of labeled data decreases. However, when the unlabeled data is used together, the performance gap is not significant, which means that unlabeled data plays an important role for classification when the number of labeled data is extremely limited. In some datasets, such as CIFAR-10, the classification performance is not good. This is because the extracted features were not well grouped by class (Figure 1-c). When more appropriate feature extraction methods are used together, better results can be obtained.

5.2. Semi-supervised Continuous Learning

The results of semi-supervised continuous learning using the NSynth dataset are shown in Figure 2. When unlabeled data is used together with the LPART, the classification performance per epoch rapidly increases to 90% and converges without catastrophic forgetting (Figure 2-a). However, when only the labeled data was used, the classification performance increased slowly and shows a similar tendency to FAM's. It

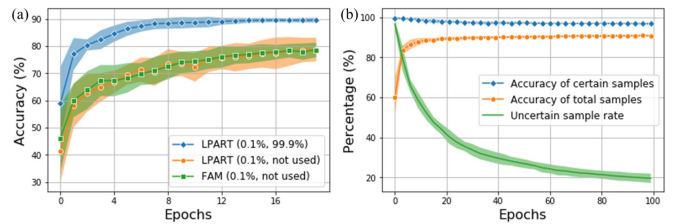


Fig. 2. (a) Semi-supervised continuous learning result by epochs. (b) Uncertain sample rate and classification accuracies using the uncertainty measurement methods by epochs. The probability of the labeled sample was set to 0.1%.

confirms that the proposed model for semi-supervised continuous learning works properly.

We also filtered out the uncertain classification results using the uncertainty measurement methods described in Section 3.4. Thresholds for two uncertainty scores were set appropriately, and only the results with the scores below these thresholds remained (Figure 2-b). The reliable results always show high performance and the number of uncertain samples continue to decrease as the learning progresses. This method is useful for applications where classification errors are fatal, and it allows selective use of reliable results in situations where labeled data is scarce and its collection is difficult.

6. CONCLUSIONS

In the present study, we proposed a novel approach for semi-supervised continuous learning based on ART networks. We applied the label propagation mechanism to the ART network and evaluated it with various datasets and experimental settings to demonstrate its effectiveness. Uncertainty measures can also be used to filter out unreliable classification results. The limitation of this study is that a pre-trained feature extractor should be used, and the quality of the extracted feature can affect the overall performance. In the future work, we will incorporate an end-to-end training of the entire system by applying a continuous learning method to the feature extractor.

7. REFERENCES

- [1] Xiaojin Zhu and Andrew B. Goldberg, "Introduction to semi-supervised learning," *Synthesis lectures on artificial intelligence and machine learning*, vol. 3, no. 1, pp. 1–130, 2009.
- [2] Michael McCloskey and Neal J Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in *Psychology of learning and motivation*, vol. 24, pp. 109–165. Elsevier, 1989.
- [3] Stephen Grossberg, "Competitive learning: From interactive activation to adaptive resonance," *Cognitive science*, vol. 11, no. 1, pp. 23–63, 1987.
- [4] Ahmet Iscen, Giorgos Toliass, Yannis Avrithis, and Ondrej Chum, "Label propagation for deep semi-supervised learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5070–5079.
- [5] Matthijs Douze, Arthur Szlam, Bharath Hariharan, and Hervé Jégou, "Low-shot learning with large-scale diffusion," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3349–3358.
- [6] Diederik P. Kingma and Max Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [7] Gail A. Carpenter, Stephen Grossberg, and David B. Rosen, "Fuzzy art: Fast stable learning and categorization of analog patterns by an adaptive resonance system," *Neural networks*, vol. 4, no. 6, pp. 759–771, 1991.
- [8] Gail A. Carpenter, Stephen Grossberg, and John H. Reynolds, "Artmap: Supervised real-time learning and classification of nonstationary data by a self-organizing neural network," *Neural networks*, vol. 4, no. 5, pp. 565–588, 1991.
- [9] Gail A. Carpenter, Stephen Grossberg, Natalya Markuzon, John H. Reynolds, David B. Rosen, et al., "Fuzzy artmap: A neural network architecture for incremental supervised learning of analog multidimensional maps," *IEEE Transactions on neural networks*, vol. 3, no. 5, pp. 698–713, 1992.
- [10] Xiao liang Tang and Min Han, "Semi-supervised bayesian artmap," *Appl Intell*, vol. 33, pp. 302–317, 2010.
- [11] G. C. Anagnostopoulos, M. Georgiopoulos, S. J. Verzi, and G. L. Heileman, "Reducing generalization error and category proliferation in ellipsoid artmap via tunable misclassification error tolerance: boosted ellipsoid artmap," in *Proceedings of the 2002 International Joint Conference on Neural Networks. IJCNN'02 (Cat. No.02CH37290)*, May 2002, vol. 3, pp. 2650–2655 vol.3.
- [12] Boaz Vigdor and Boaz Lerner, "The bayesian artmap," *IEEE Transactions on Neural Networks*, vol. 18, no. 6, pp. 1628–1644, 2007.
- [13] Xianxu Hou, Linlin Shen, Ke Sun, and Guoping Qiu, "Deep feature consistent variational autoencoder," in *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2017, pp. 1133–1141.
- [14] Yunchen Pu, Zhe Gan, Ricardo Henao, Xin Yuan, Chunyuan Li, Andrew Stevens, and Lawrence Carin, "Variational autoencoder for deep learning of images, labels and captions," in *Advances in neural information processing systems*, 2016, pp. 2352–2360.
- [15] Theofanis Karaletsos, Serge Belongie, and Gunnar Rätsch, "Bayesian representation learning with oracle constraints," *arXiv preprint arXiv:1506.05011*, 2015.
- [16] Xiaojin Zhu and Zoubin Ghahramani, "Learning from labeled and unlabeled data with label propagation," Tech. Rep., School of Computer Science, Carnegie Mellon University, 2002.
- [17] George J Klir and Mark J Wierman, *Uncertainty-based information: elements of generalized information theory*, vol. 15, Physica, 2013.
- [18] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al., "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [19] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y. Ng, "Reading digits in natural images with unsupervised feature learning," 2011.
- [20] Alex Krizhevsky, Geoffrey Hinton, et al., "Learning multiple layers of features from tiny images," Tech. Rep., Citeseer, 2009.
- [21] Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Mohammad Norouzi, Douglas Eck, and Karen Simonyan, "Neural audio synthesis of musical notes with wavenet autoencoders," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 1068–1077.
- [22] Laurens van der Maaten and Geoffrey Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.