



Large-Margin Representation Learning for Texture Classification

Jonathan de Matos^{a,d,**}, Luiz Eduardo Soares Oliveira^c, Alceu de Souza Britto Junior^{b,d}, Alessandro Lameiras Koerich^a

^aÉcole de Technologie Supérieure (ÉTS), Université du Québec, 1100 Notre-Dame Street West, Montreal, QC, H3C 1K3, Canada

^bPontifícia Universidade Católica do Paraná (PUCPR), Rua Imaculada Conceição, 1155, Curitiba, PR, 80215-901, Brazil

^cUniversidade Federal do Paraná (UFPR), Rua Coronel Francisco Heraclito dos Santos, 100, Curitiba, PR, 81530-900, Brazil

^dUniversidade Estadual de Ponta Grossa (UEPG), Avenida General Carlos Cavalcanti, 4748, Ponta Grossa, PR, 84030-900, Brazil

Article history:

Fully Convolutional Networks; Large-margin Classifier; Feature Extraction; Texture

ABSTRACT

This paper presents a novel approach combining convolutional layers (CLs) and large-margin metric learning for training supervised models on small datasets for texture classification. The core of such an approach is a loss function that computes the distances between instances of interest and support vectors. The objective is to update the weights of CLs iteratively to learn a representation with a large margin between classes. Each iteration results in a large-margin discriminant model represented by support vectors based on such a representation. The advantage of the proposed approach w.r.t. convolutional neural networks (CNNs) is two-fold. First, it allows representation learning with a small amount of data due to the reduced number of parameters compared to an equivalent CNN. Second, it has a low training cost since the backpropagation considers only support vectors. The experimental results on texture and histopathologic image datasets have shown that the proposed approach achieves competitive accuracy with lower computational cost and faster convergence when compared to equivalent CNNs.

© 2022 Elsevier Ltd. All rights reserved.

1. Introduction

Convolutional neural networks (CNNs) have been established as the state-of-the-art approach to computer vision. CNNs achieve high accuracy due to how the convolutional layers filter images, the number of parameters available to learn representations [1][2], and the large datasets [3] used in their training. They usually present high accuracy in object recognition problems, e.g., cars, people, animals, numbers. Networks pretrained with large datasets of objects can be reused in other contexts. There are two ways of transferring networks between contexts, fine-tuning or using them as feature extractors. The fine-tuning procedure allows fast training with fewer data because pre-trained filters are already able to identify some patterns. Although it reduces training effort, data may be insufficient even for fine-tuning in small datasets. Furthermore, when

working with datasets where textural information prevails to the detriment of shape and spatial characteristics, the first layers of pretrained CNNs may not respond well to new patterns. There are two additional problems: (i) deeper layers also have more problem-specific knowledge; (ii) there is the need for image size adaptation, which can slow the training [4]. Using CNNs as feature extractors, there is no filter update. Instead, the activation maps of a specific layer are used as a feature vector for images of a new context, and an alternative classifier is trained on them. They work similarly to a handcrafted feature extractor and neither adapt to the new patterns nor generate features targeted to facilitate the classification task. Small and simple models like these are more suitable for classifying non-object and small datasets since they do not require data for representation learning but only for training a discriminant.

This paper proposes a novel large-margin representation learning that overcomes most of the problems mentioned above related to CNNs and handcrafted feature extractors. For this purpose, it addresses the following questions: a) is it possible to learn representations and discriminants on small-size texture datasets? b) is it possible to speed up the training convergence

**Corresponding author

e-mail: jonathan@uepg.br (Jonathan de Matos),
luiz.oliveira@ufpr.br (Luiz Eduardo Soares Oliveira),
alceu@ppgia.pucpr.br (Alceu de Souza Britto Junior),
alessandro.koerich@etsmt1.ca (Alessandro Lameiras Koerich)

while achieving high accuracy? The proposed approach uses a short sequence of convolutional layers (CLs) to learn representation for texture classification. The CLs feed a large-margin discriminant that, in turn, provides information to update the CLs' weights and increase the decision margin. A novel loss function calculates the distance between instances in the decision frontier and anchors. The backpropagation algorithm minimizes the loss function while enlarging the margin between classes. The novelty of our approach is that it employs only instances that violate the decision margin to train the CLs and produce latent representations. That speeds up the convergence of the backpropagation algorithm to suitable latent representation and discriminant. In addition, it uses fewer parameters than a conventional CNN, allowing training in small datasets, and it performs well on non-object context recognition. Otherwise, the CNNs would still be more effective. The proposed approach was evaluated on a synthetic dataset based on Gaussian distributions, texture datasets, and three histopathological image (HI) datasets.

The main contributions of this paper are: (i) an approach that trains CLs from scratch with little data; (ii) A representation learning method that adapts itself to the characteristics of different texture datasets; (iii) a computational efficient training technique that uses only support vectors (SVs) in each iteration instead of all training instances; (iv) A fast convergence method compared to methods used for training conventional CNNs; (v) resilience to imbalanced data; (vi) A detailed comparison of the performance achieved by the proposed approach with approaches using handcrafted features and CNNs.

This paper is organized as follows. Section 2 presents some fundamental concepts. Section 3 describes the proposed approach and the experiment setup. In Section 4, we present some experimental results and discuss the advantages of our method. The conclusions and the main contributions of our approach are presented in the last section.

2. Related Works

The first CLs of CNNs are suitable to identify general and straightforward patterns such as textures. Their training help to make them adaptable to motifs of each context, so it is worthwhile to train these first CLs and use only them in more simple image contexts. One may find different contributions in the literature to adapt pre-trained CNN models to a new domain.

Cimpoi et al. [5] proposed the Fisher vector CNN (FV-CNN), which pools the last CL of a pretrained network, using it as a feature vector. The pooling allows using input images without resizing them to fit a fully connected CNN (FC-CNN). They used the FV-CNN features as input to an SVM and compared the FC-CNN and SIFT features. The FV-CNN showed to be a good texture descriptor, performing well on several benchmarks. The texture CNN (TCNN)[4] uses a similar approach to the FV-CNN, but with the classification accomplished by a sequence of FC layers. Its difference to a traditional FC-CNN is the energy layer, a global average pooling (GAP) at the last CL whose outputs go to the FC layers. The SVM replacement by the FC layers as a classifier, in contrast to FV-CNN, permits training the CLs in conjunction with the classifier. This last approach does not require a pretrained CNN as it can train it with

the FC layers, but it can explore transfer learning on low data scenarios, such as some types of medical images acquired by MRI, CT, radiography, or microscopy that have textural characteristics [6]. Microscopy images allows the observation of tissues, their cells, and their nuclei, with a more texture oriented appearance. They are a challenging classification problem due to differences on staining process and small number of samples. Therefore, instead of using large CNNs, one can use only their first layers that are more texture-aware as an FCN, intercepting the output and treating it as a feature vector. This latent representation can be used with a simpler classifier with fewer parameters than FC layers of a CNN.

One can also use the deep metric learning (DML) concepts to train the FCNs to produce better filters and a representation more suitable to the new classifier.

Metric learning is an alternative to classification methods in situations where the number of classes is high, in the order of thousands, or the number of samples of each category is low [7]. Two examples of this situation are face recognition and signature verification. The idea of the dissimilarity metric learning is to learn a representation to identify two samples' similarities, usually employing a distance metric, e.g., Euclidean distance. DML approaches allow metric and representation learning simultaneously. Learning representation is possible with deep learning due to the training of parameters used in its layers, mostly the convolutional ones. Chopra et al. [7] presented a method that uses siamese convolutional neural networks and the energy-based model. It consists of using raw energy values instead of probabilistic normalized ones. Their concept of energy is analogous to the one of Andrearczyk and Whelan [4]. The advantage of using CNNs is that they provide end-to-end training that learns low and high-level features and results in shift-invariant detectors. Their method aggregates the energy output of each siamese network into one neural network trained with contrastive loss. The loss allows training the system together (siamese networks and the single neural network), improving the sample representation on the energy layer for texture classification.

DML algorithms can be split into pair-based or proxy-based. The former, like contrastive loss, aims to minimize intra-class distance and maximize inter-class distance. However, this approach has prohibitive computational complexity and pairs that do not contribute to the training. Therefore, several works have addressed these issues [8]. The ranked list loss method [9] relies on a threshold margin that gives more attention, using weighting, to the samples that maximize the margins between opposite classes. The margin determines the negative points that are too close to the query and violates most of the margin. These are the negative selected points. On the other hand, proxy-based methods [10] [11] create an instance representing a set of instances of the same class. It reduces the number of training instances and avoids noise and outliers.

3. Large-Margin Fully Convolutional Network (LMFCN)

The proposed approach is pair-based, and it selects the most effective instances for the training procedure, reducing the training complexity and discarding irrelevant instances. It also

selects specific examples as anchors to calculate a novel loss function, which speeds up training. Although the proposed approach shares ideas of DML, our application context is different, with more instances per class and fewer classes than the usual DML context.

The proposed method has three components: a fully convolutional network (FCN) with a global average pooling (GAP), a large-margin classifier, and a novel loss function, made up of three terms. The LMFCN¹ uses a sequence of CL acting as a filter bank to learn representation from data used as input to a large-margin classifier. It acts as an end-to-end image classifier, like a CNN but requires less training data to learn high discriminant representations. Its advantage is to enable filter training and make the latent representation more suitable for the classifier. In addition, the backpropagation algorithm trains the filters with a particular loss function that uses the distance between instances of interest and their anchors provided by the large margin classifier to improve representation learning.

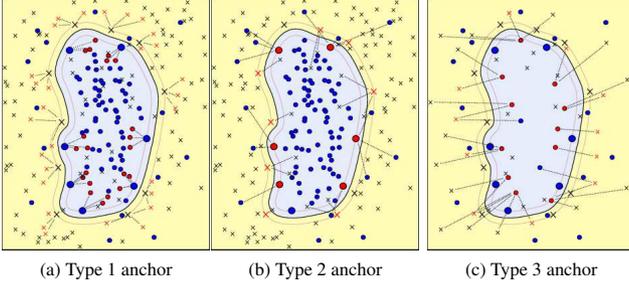


Fig. 1: A two-class latent space with decision boundary computed by an RBF SVM. Circles are samples of class 0 and crosses, class 1. Bigger symbols mean the SVs for each class. In blue is the region of class 0, and in yellow is class 1. Dashed straight black lines link the instances to their anchors. In red are the anchors for three different situations on the calculations. (a) Reference anchors to the SVs; (b) Anchors used to move the misclassified instances; (c) Anchors used to increase the separation of instances from opposite classes.

3.1. Training Procedure and Loss Function

The learning algorithm uses the concept of anchors to guide the training. Anchors are instances from the training set used as references by the loss function to calculate the distance to instances of interest. The backpropagation algorithm minimizes such distances during training. Fig. 1 presents a latent space split into two regions by an RBF SVM classifier and the three types of anchors.

Type 1 anchors (red circles and crosses) are the correctly classified instances closest to the support vectors (SVs). The LMFCN attempts to maximize the margin by pushing the SVs in the direction of such anchors. The learning algorithm minimizes the distance between them and the SVs. In Fig. 1a, the dotted straight lines linking SVs identify the three anchors of each SV and help visualize the effect of the distance minimization.

Type 2 anchors (Fig. 1b) are the SVs closest to misclassified examples used to move such instances in direction to the right side of the decision boundary.

Type 3 anchors are the closest correctly classified instances of the opposite class, as shown in Fig. 1c. Type 3 anchors help

to maximize the distance between samples of different classes. The number of anchors, regardless of their type, is a hyperparameter.

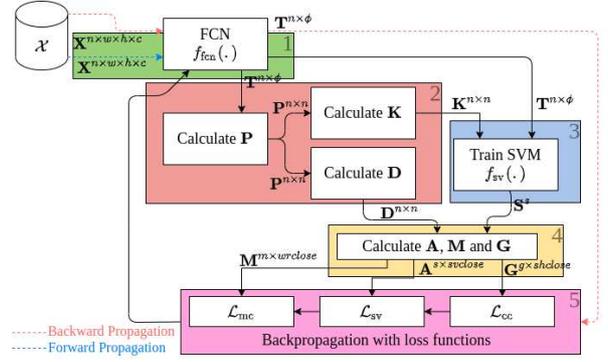


Fig. 2: An overview of the training scheme of the LMFCN.

The training procedure starts with a dataset of size n denoted as $\mathcal{X} = \{\mathbf{X}^t, o^t\}_{t=1}^n$, where t indexes images \mathbf{X}^t of width w , height h and c channels in \mathcal{X} , and $o^t \in \mathbb{N} : [0, 1]$ is its expected output. All images from \mathcal{X} are fed to CLs denoted as $f_{\text{fcn}}(\cdot)$ and produces a matrix $\mathbf{T}^{n \times \phi}$, with ϕ being the size of the latent representation (step 1 at Fig. 2). The algorithm uses matrix \mathbf{T} to calculate matrix \mathbf{P} (Eq. (1)), which in turn is used to calculate matrices $\mathbf{K}^{n \times n}$ and $\mathbf{D}^{n \times n}$ using Eqs. (2) and (3), respectively (step 2 at Fig. 2). \mathbf{K} is a RBF kernel matrix used to train a large margin classifier $f_{\text{sv}}(\cdot)$, which produces the output y^t for each element t of \mathcal{X} , and also provides the set of support vectors (SVs) indexes $\mathcal{S} = \{s^u\}_{u=1}^v$, where $s^u \subset \mathbb{N} : [0, n[$ and v is the number of SVs (step 3 at Fig. 2). \mathbf{D} is essential to the rest of the algorithm as it contains the pairwise distance of all instances and is used to create the anchor matrices.

$$p_{ij} = \sum_{b=0}^{\phi} (T_{ib} - T_{jb})^2 \quad (1)$$

$$k_{ij} = \exp(-\gamma p_{ij}) \quad (2)$$

$$d_{ij} = \sqrt{p_{ij}} \quad (3)$$

Type 1 anchors (step 4 at Fig. 2) are obtained from matrix \mathbf{D} , the set \mathcal{S} , the expected output o^t , and the output from $f_{\text{sv}}(f_{\text{fcn}}(\mathbf{X}^t))$, as shown in Eq. (4):

$$e_{ij} = \begin{cases} \tau & \text{if } i = j \text{ or } j \in \mathcal{S} \text{ or } o^i \neq o^j \text{ or } o^j \neq f_{\text{sv}}(f_{\text{fcn}}(\mathbf{X}^j)) \\ d_{ij} & \text{otherwise} \end{cases} \quad (4)$$

where e_{ij} is either the distance between SVs and their anchors (d_{ij}) or a large constant τ that indicates that there is no SV-anchor relation (uninteresting instances), $f_{\text{fcn}}(\mathbf{X}^t)$ denotes the latent representation generated by the FCN for an input \mathbf{X}^t , $f_{\text{sv}}(\cdot)$ is the output predicted by the large-margin classifier.

Furthermore, we also define a sorting function $f_{\text{argsort}}(\cdot)$ that takes a vector as input and returns a vector of indices sorted in increasing order. So, the anchor matrix $\mathbf{A}^{s \times n}$ is calculated using Eq. (5).

$$\mathbf{a}_u = f_{\text{argsort}}(\mathbf{e}_{s^u}) \text{ with } s^u \in \mathcal{S} \text{ and } u \in [0, |\mathcal{S}[\quad (5)$$

Likewise \mathbf{A} , we calculate with the Eqs. (7) and (9) the matrices $\mathbf{M}^{m \times n}$ and $\mathbf{G}^{g \times n}$ for type 2 and 3 anchors, respectively (step 4 at Fig. 2). We define a set of indexes to the misclassified

¹<https://github.com/jonathandematos/lmfcn>

and correct classified instances from \mathcal{X} as $\mathbf{Q} = \{q \mid q \in \mathbb{N} : [0, n[\wedge f_{sv}(f_{fcn}(\mathbf{X}^q)) \neq o^q]\}$, and $\mathbf{R} = \{r \mid r \in \mathbb{N} : [0, n[\wedge r \notin S \cup \mathbf{Q}]\}$, respectively.

$$z_{ij} = \begin{cases} d_{ij} & \text{if } i \in \mathbf{Q} \text{ and } j \in \mathbf{S} \text{ and } i \neq j \\ \tau & \text{otherwise} \end{cases} \quad (6)$$

$$\mathbf{m}_i = f_{\text{argsort}}(\mathbf{z}_{q^i}) \text{ with } q^i \in \mathbf{Q} \text{ and } i \in [0, |\mathbf{Q}|[\quad (7)$$

$$h_{ij} = \begin{cases} d_{ij} & \text{if } i \in \mathbf{R} \text{ and } j \in \mathbf{R} \text{ and } o^i \neq o^j \\ \tau & \text{otherwise} \end{cases} \quad (8)$$

$$\mathbf{g}_i = f_{\text{argsort}}(\mathbf{h}_{r^i}) \text{ with } r^i \in \mathbf{R} \text{ and } i \in [0, |\mathbf{R}|[\quad (9)$$

where z_{ij} and h_{ij} in Eqs. (6) and (8) are the distance between misclassified instances and the SVs (d_{ij}) and the distance between correct classified instances from opposite classes (d_{ij}), respectively. Again, τ is a large constant that indicates uninteresting instances.

Step 5 of the algorithm is the backpropagation of images from \mathcal{X}^S , \mathcal{X}^Q and \mathcal{X}^R through the FCN to calculate the gradients and use them with the loss functions to update the CL weights. Steps 1 to 5 define one epoch. After step 5, the process restarts from step 1, so the latent representation is computed again, as though the matrices \mathbf{P} , \mathbf{K} and \mathbf{D} and the large-margin discriminant is retrained on such a updated latent representation, producing another set of SVs allowing recalculation of matrices \mathbf{A} , \mathbf{M} and \mathbf{G} . Therefore, the elements of matrix \mathbf{T} are different from the previous epoch because of the updated FCN weights.

The proposed loss function relies on the similarity between examples, and it has three terms, as shown in Eq. (10). It aims at finding a latent representation that maximizes the margin (\mathcal{L}_{sv}) while pushing misclassified examples towards the right side of the decision boundary (\mathcal{L}_{mc}) and moving well-classified examples farther away from the decision boundary (\mathcal{L}_{cc}).

$$\mathcal{L} = \mathcal{L}_{sv} + \mathcal{L}_{mc} + \mathcal{L}_{cc} \quad (10)$$

\mathcal{L}_{sv} calculates the sum of distances between SVs and its anchors using matrix \mathbf{A} , as shown in Eq. (11). As a consequence, the gradients are affected only by the instances which are SVs, not by the type 1 anchors, as they were already generated and stored in \mathbf{T} . The backpropagation procedure updates the weights in a way that the latent representation generated by $f_{fcn}(\cdot)$ has the smallest possible distance to the fixed values of type 1 anchors at the current epoch. Therefore, the weights are updated to move the SVs and not the anchors.

$$\mathcal{L}_{sv} = \frac{\sum_{i=0}^{|\mathbf{S}|} \sum_{j=1}^{sv_{\text{close}}} [f_{fcn}(\mathbf{X}^{s^i}) - \mathbf{t}_{a_{ij}}]^2}{|\mathbf{S}|} \quad (11)$$

where sv_{close} is the number of anchors to use for each SV, $f_{fcn}(\cdot)$ is the FCN updated by the backpropagation, \mathbf{X}^{s^i} is the matrix that represents the s^i image from \mathcal{X} , a_{ij} is an index pointing to an anchor instance in the input set \mathcal{X} , $\mathbf{t}_{a_{ij}}$ is the latent representation of an image $\mathbf{X}^{a_{ij}}$, and $|\mathbf{S}|$ is the number of SVs.

The training algorithm computes the loss over the entire set of SVs as a single batch. It is also possible to use mini-batches, but considering small-size datasets and an FCN with compact architecture that yields a low-dimensional latent representation,

this is unnecessary. In the next epoch, the updated weights will affect the generation of the latent representation of the entire dataset. Therefore, the latent representation of anchors also changes, and the training algorithm builds a new set of anchors.

\mathcal{L}_{mc} calculates the summation of distances between misclassified instances and their anchors using matrix \mathbf{M} , as shown in Eq. (12). We also use all misclassified instances as a single batch, although there is no limitation to performing it in mini-batches. The influence of \mathcal{L}_{mc} in the training algorithm is the weight updating that minimizes the distance between the misclassified instances and their closest SVs. Consequently, the misclassified instances are pushed towards the right side of the decision boundary.

$$\mathcal{L}_{mc} = \frac{\sum_{i=0}^{|\mathbf{Q}|} \sum_{j=0}^{wr_{\text{close}}} [f_{fcn}(\mathbf{X}^{q^i}) - \mathbf{t}_{m_{ij}}]^2}{|\mathbf{Q}|} \quad (12)$$

where wr_{close} is the number of anchors for each misclassified instance, $f_{fcn}(\cdot)$ is the FCN updated by the backpropagation, \mathbf{X}^{q^i} is the matrix that represents the q^i image from \mathcal{X} , m_{ij} is an index pointing to an anchor instance in the input set \mathcal{X} , $\mathbf{t}_{m_{ij}}$ is the latent representation of an image $\mathbf{X}^{m_{ij}}$, and $|\mathbf{Q}|$ is the number of misclassified instances.

When looking at only a single misclassified instance, the weight updating may not be enough to move such an instance to the right side of the decision boundary because its anchors are SVs right at the margin of the decision boundary. Despite that, the misclassified example can become an SV in the following training epoch.

\mathcal{L}_{cc} represents the distance between well-classified instances and their anchors. Since we want to maximize such a distance, Eq. (13) calculates the inverse of such a distance and incorporates it in the loss function, which is minimized during training.

$$\mathcal{L}_{cc} = \frac{|\mathbf{R}|}{\sum_{i=0}^{|\mathbf{R}|} \sum_{j=0}^{sh_{\text{close}}} [f_{fcn}(\mathbf{X}^{r^i}) - \mathbf{t}_{g_{ij}}]^2} \quad (13)$$

where sh_{close} is the number of anchors for each correctly classified instance, $f_{fcn}(\cdot)$ is the FCN updated by the backpropagation, \mathbf{X}^{r^i} is the matrix that represents the r^i image from \mathcal{X} , g_{ij} is an index pointing to an anchor instance in the input set \mathcal{X} , $\mathbf{t}_{g_{ij}}$ is the latent representation of an image $\mathbf{X}^{g_{ij}}$, and $|\mathbf{R}|$ is the number of correct classified instances.

The number of anchors used for each training instance is controlled by sv_{close} , wr_{close} , and sh_{close} in Eqs. (11) (12), and (13), respectively. Usually, complex classification problems require a higher number of anchors. Furthermore, using all three terms of the proposed loss function in the training process may not always be necessary. \mathcal{L}_{sv} alone already leads to good representations as such a loss is directly related to SVs and margin maximization. The computational cost for computing this term of the loss function is not high and decreases as the number of instances used by the backpropagation algorithm reduces at each training epoch. Computing \mathcal{L}_{mc} is not also expensive because the number of incorrect classified examples tends to decrease over the training epochs. On the other hand, computing

\mathcal{L}_{cc} can become very expensive because the number of well-classified instances tends to increase over the training epochs. Therefore, the term \mathcal{L}_{cc} should be used wisely, preferably only on challenging problems where just the other two terms of the proposed loss function may not lead to a low training error.

3.2. FCN Architecture

The FCN is a sequence of CLs similar to ones used in CNNs, used as filter banks to learn representation related to textures, as presented in Table 1. Pooling layers follow these layers to reduce the input size progressively. The deeper the layers, the narrower is the latent representation, but with an increasing number of channels, which allows more filters combination, increasing the complexity of the representation. At the end of the CLs, a GAP layer builds a latent representation where each element represents the response of a filter combination to a texture, measuring how much it happened and not its position on the image. The GAP layer also makes the output dimension independent of the input size, and the latent representation will always have the same number of channels at the end of the FCN. Such a latent representation feeds a large-margin discriminant to learn classification tasks. We compared our approach with two CNNs with identical architecture (Table 1), but with a sequence of fully connected (FC) layers as discriminant after the GAP layer. We employed binary cross-entropy (BCE) loss and Hinge loss for binary problems and cross-entropy loss in multi-class.

Table 1: FCN used on the LMFCN and in the CNN comparison. (w : width, h : height, c : number of channels, ϕ : dimension of the latent representation.)

Layer	Input	Output
Convolutional Layer	$w \times h \times c$	$w \times h \times 64$
Max Pooling	$w \times h \times 64$	$w/2 \times h/2 \times 64$
Convolutional Layer	$w/2 \times h/2 \times 64$	$w/2 \times h/2 \times 128$
Max Pooling	$w/2 \times h/2 \times 128$	$w/4 \times h/4 \times 128$
Convolutional Layer	$w/4 \times h/4 \times 128$	$w/4 \times h/4 \times \phi$
Global Average Pooling	$w/4 \times h/4 \times \phi$	$1 \times 1 \times \phi$

Batch Normalization and ReLU after each CL.

The large margin classifier of the LMFCN is a support vector machine (SVM) with an RBF kernel, which Gram matrix \mathbf{K} is obtained by Eq. (2) calculated jointly with the distance matrix \mathbf{D} . Part of the classifier calculation is reused and performed in GPU. We chose the precomputed RBF kernel due to its ease of computation using the GPU resources and space separation capacity. Although a linear kernel has reduced computational cost, it would require more training of the FCN weights to provide features with more class separation. In the preliminary studies, we compared the two kernel approaches and verified the advantage of the RBF.

A large-margin discriminant is inherently binary, and to deal with multiclass problems, we adopted the one-vs-all (OVA) approach, which reduces multiclass problems into multiple binary classification problems. In the training stage, one provides all instances to all n_c pairs of LMFCNs. After training all the n_c LMFCNs, we discard the discriminants, keeping only the FCNs. Then, we train a new multiclass SVM. The new classifier is trained using a latent representation with $\phi \times n_c$, which is the concatenation of all FCN latent representations. At the end of the training process, the full model comprises n_c FCNs with an ϕ -dimensional output and a multiclass SVM with a latent representation of $\phi \times n_c$ dimensions. The multiclass SVM holds

internally multiple binary SVMs with RBF kernel on OVA configuration. Although this approach seems similar to using multiple SVMs trained with the FCNs, the new classifiers have access to a latent representation that is better fitted for them.

The computational cost for training the LMFCN and equivalent CNNs up to the GAP layer is proportional to the number and resolution of input images ($n \times w \times h \times c$) and the number weights (α_{fcn}). The LMFCN replaces the FC layers (computational cost of $\alpha_{fcn} \times n$), by an SVM, which requires kernel calculation (n^2), and sequential minimal optimization (SMO), which requires n^3 in the worst case. Assuming a small training set, the SMO cost is lower than $n \times \alpha_{fcn}$. However, for large datasets, where n^2 is greater than α_{fcn} , the problem becomes more suitable for conventional CNN architectures.

The main advantage of the LMFCN is using only the SVs in the backpropagation, which reduces its computational cost from n to $|\mathcal{S}|$. The CLs used on both LMFCN and conventional CNNs have several trainable parameters. Therefore, having only SVs as training instances reduces the number of training instances and speeds up the training, making the loss function converge faster within a few epochs.

4. Experimental Results and Discussion

The evaluation of the proposed LMFCN is carried out on images with texture characteristics and low training data availability. We used a synthetic dataset of images generated from two Gaussian distributions with striped patterns, misc and fabric categories of the Salzburg Texture Image Database [12], BreakHis [13], BACH [14] and CRC [15] datasets, which comprise histopathological images (HIs).

Fig. 3 shows a training graph of the Gaussian images dataset over 20 epochs with the LMFCN. The best-balanced accuracy values for train, validation, and test correspond to the peak validation accuracy. It is possible to observe that the value of losses and the number of SVs goes down over the epochs, and the accuracy rises. That shows that the SVs are getting closer to their instances due to the loss reduction. In addition, the decrease in the number of SVs indicates that a better separability is being achieved.

Table 2 compares the results achieved by the LMFCN with other CNN architectures on five datasets. Two CNNs have an architecture similar to the LMFCN, but they use the Hinge loss (CNN-H) and the binary cross-entropy loss (CNN-CE). These two CNNs are equivalent to TCNNs [4] since they use a short sequence of CLs and a GAP to produce energy values. We also used pretrained ResNet18 and InceptionV3 as feature extractors. We limited the number of training epochs to 100 and 15 epochs for the CNNs and LMFCN. The discriminant layers of the ResNet18 and the InceptionV3 have the same parameters as the LMFCN. Furthermore, we used PCA to reduce the feature vectors generated by these two CNNs from 512 and 2048 to 16, the same size used by the other architectures. Overall, the LMFCN converges to a latent representation that generalizes well much faster (less than ten epochs for all datasets) than different CNN architectures. Furthermore, in 3 out of 5 datasets, the accuracy achieved by the LMFCN on the test sets is higher than that achieved by other CNNs. For the two other datasets, the difference in accuracy is almost negligible.

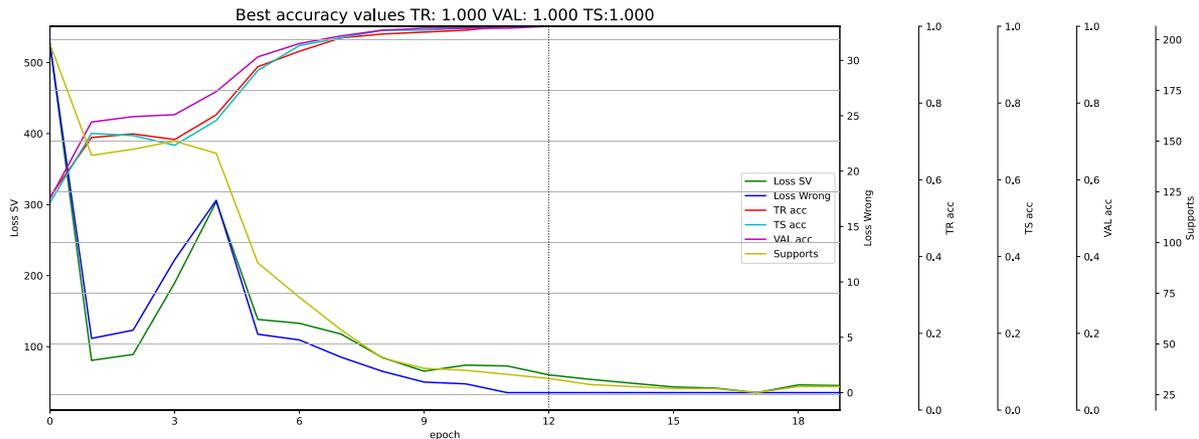


Fig. 3: Training of the LMFCN using synthetic Gaussian images with a 2-dimensional latent representation. Parameters: $sv_{close}=5$, $wr_{close}=1$, and $sh_{close}=0$.

Table 2: Balanced accuracy for the LMFCN and equivalent CNN architectures. Epoch refers to training epoch where the best validation accuracy was achieved. NA: Not Applicable, for models that are trained in one single epoch

Dataset	Architecture	Training	Validation	Test	Epoch
BACH	LMFCN	0.8811 ± 0.0426	0.7857 ± 0.0221	<u>0.8056 ± 0.0203</u>	6
	CNN-CE	0.8926 ± 0.0135	<u>0.8200 ± 0.0291</u>	0.7864 ± 0.0212	88
	CNN-H	0.8534 ± 0.0410	0.8038 ± 0.0261	0.7838 ± 0.0242	68
	ResNet18	<u>0.9981 ± 0.0026</u>	0.7596 ± 0.0122	0.7584 ± 0.0270	NA
	InceptionV3	0.6295 ± 0.0746	0.5924 ± 0.0625	0.6049 ± 0.0630	NA
BreakeHis	LMFCN	<u>0.9882 ± 0.0064</u>	<u>0.9442 ± 0.0137</u>	<u>0.8942 ± 0.0372</u>	5
	CNN-CE	0.8926 ± 0.0092	0.9122 ± 0.0226	0.8926 ± 0.0140	91
	CNN-H	0.8560 ± 0.0041	0.8856 ± 0.0166	0.8638 ± 0.0042	91
	ResNet18	0.9777 ± 0.0062	0.8034 ± 0.0120	0.8262 ± 0.0122	NA
	InceptionV3	0.6058 ± 0.0207	0.6064 ± 0.0209	0.5949 ± 0.0207	NA
CRC	LMFCN	0.9924 ± 0.0046	0.9914 ± 0.0024	0.9914 ± 0.0060	6
	CNN-CE	0.9828 ± 0.0070	<u>0.9934 ± 0.0025</u>	0.9880 ± 0.0111	31
	CNN-H	0.9928 ± 0.0038	0.9924 ± 0.0027	<u>0.9928 ± 0.0070</u>	30
	ResNet18	<u>0.9991 ± 0.0013</u>	0.9683 ± 0.0081	0.9680 ± 0.0150	NA
	InceptionV3	<u>0.7718 ± 0.0352</u>	0.7670 ± 0.0375	0.7842 ± 0.0355	NA
Salzburg	LMFCN	0.9842 ± 0.0049	<u>0.9446 ± 0.0083</u>	<u>0.9230 ± 0.0081</u>	8
	CNN-CE	0.8966 ± 0.0103	0.8762 ± 0.0100	0.8602 ± 0.0115	91
	CNN-H	0.8258 ± 0.0257	0.8292 ± 0.0277	0.8046 ± 0.0314	84
	ResNet18	<u>0.9946 ± 0.0033</u>	0.9046 ± 0.0221	0.8932 ± 0.0088	NA
	InceptionV3	<u>0.6680 ± 0.0332</u>	<u>0.6607 ± 0.0323</u>	<u>0.6644 ± 0.0473</u>	NA
Gaussian	LMFCN	<u>1.0000 ± 0.0000</u>	<u>1.0000 ± 0.0000</u>	0.9990 ± 0.0022	2
	CNN-CE	0.9950 ± 0.0061	<u>1.0000 ± 0.0000</u>	<u>1.0000 ± 0.0000</u>	72
	CNN-H	<u>1.0000 ± 0.0000</u>	<u>1.0000 ± 0.0000</u>	<u>1.0000 ± 0.0000</u>	56
	ResNet18	0.9740 ± 0.0109	0.6126 ± 0.0198	0.6176 ± 0.0257	NA
	InceptionV3	0.5102 ± 0.0235	0.5311 ± 0.0185	0.5021 ± 0.0258	NA

Table 3 shows average results of balanced accuracy of five repetitions comparing the results achieved by the LMFCN with shallow approaches that employ an SVM and three handcrafted feature extractors: Local Binary Pattern (LBP), Gray Level Co-occurrence Matrix (GLCM), and Parameter Free Threshold Adjacency Statistics (PFTAS) [13].

For a fair comparison, the LMFCN uses a 59-dimensional latent representation, the same dimension of the smallest feature vector, obtained with the LBP with uniform patterns. PFTAS and GLCM produce 162- and 169-dimensional feature vectors, respectively. PFTAS achieved the best accuracy among the handcrafted feature extractors on three datasets (BACH, BreakeHis, and CRC). LBP achieved the best accuracy on Salzburg and GLCM on the Gaussian dataset. However, the LMFCN with a 59-dimensional latent representation achieves an accuracy higher than all shallow methods on all datasets, indicating its adaptability to different problems and datasets.

The multiclass experiments were carried out on three HI datasets. Our experiments also included comparisons against

Table 3: Balanced accuracy for LMFCN, GLCM, LBP and PFTAS. The best results on each subset are underlined.

Dataset	Method	Training	Validation	Test
BACH	LMFCN-59	0.9664 ± 0.0751	0.7888 ± 0.0577	0.7684 ± 0.0435
	LMFCN-162	<u>1.0000 ± 0.000</u>	<u>0.8068 ± 0.0600</u>	<u>0.7934 ± 0.0536</u>
	GLCM	0.7346 ± 0.0199	0.6309 ± 0.0346	0.6618 ± 0.0711
	LBP	0.9146 ± 0.0154	0.7455 ± 0.0442	0.7005 ± 0.0228
	PFTAS	<u>0.9899 ± 0.0037</u>	<u>0.7372 ± 0.0634</u>	<u>0.7608 ± 0.0361</u>
BreakeHis	LMFCN-59	<u>1.0000 ± 0.0000</u>	0.9639 ± 0.0247	0.9476 ± 0.0075
	LMFCN-162	<u>1.0000 ± 0.0000</u>	<u>0.9666 ± 0.0183</u>	<u>0.9595 ± 0.0050</u>
	GLCM	0.8128 ± 0.0086	0.8117 ± 0.0375	0.8076 ± 0.0071
	LBP	0.9292 ± 0.0059	0.7998 ± 0.0099	0.7925 ± 0.0098
	PFTAS	<u>0.9801 ± 0.0034</u>	<u>0.9237 ± 0.0236</u>	<u>0.9145 ± 0.0168</u>
CRC	LMFCN-59	0.9977 ± 0.0022	<u>0.9937 ± 0.0023</u>	<u>0.9883 ± 0.0084</u>
	GLCM	0.9864 ± 0.0036	0.9828 ± 0.0055	0.9853 ± 0.0101
	LBP	0.9974 ± 0.0012	0.9431 ± 0.0099	0.9479 ± 0.0061
	PFTAS	<u>1.0000 ± 0.0000</u>	<u>0.9852 ± 0.0025</u>	<u>0.9872 ± 0.0097</u>
	Salzburg	LMFCN-59	<u>0.9994 ± 0.0006</u>	<u>0.9672 ± 0.0112</u>
GLCM		0.7706 ± 0.0163	0.7257 ± 0.0186	0.7338 ± 0.0142
LBP		0.9987 ± 0.0012	0.9282 ± 0.0209	0.9116 ± 0.0208
PFTAS		0.9650 ± 0.0051	0.9096 ± 0.0124	0.8973 ± 0.0137
Gaussian		LMFCN-59	<u>1.0000 ± 0.0000</u>	<u>0.9354 ± 0.1254</u>
	GLCM	0.9583 ± 0.0028	0.9327 ± 0.0110	0.9527 ± 0.0191
	LBP	0.6303 ± 0.0091	0.6710 ± 0.0131	0.6109 ± 0.0075
	PFTAS	0.5578 ± 0.0080	0.5953 ± 0.0077	0.5608 ± 0.0076

handcrafted feature extractors GLCM, LBP, and PFTAS as though as ResNet18 and InceptionV3 as feature extractors. We also used a CNN with cross-entropy loss (CE) and similar architecture to the LMFCN but using more filters. The number of filters is proportional to the total number of parameters summing all OVA models of the LMFCN, providing a fair comparison.

Table 4 shows the average balanced accuracy for all evaluated methods. It is noticeable the superior performance of the LMFCN over all others. In the case of the CNN experiments, we allowed the training procedure to extend over 400 epochs. Although the LMFCN uses multiple models, we let each one only train for ten epochs. The metric used, balanced accuracy, helps identify problems on imbalanced datasets. Our approach presented a promising performance on BreakeHis, which has a significant difference between some classes, e.g., Ductal Carcinoma and Phyllodes Tumor. This result shows that the LMFCN performed well in imbalanced scenarios. We also noticed that the LMFCN performed well on the imbalanced OVA subproblems.

The experimental results have shown that the LMFCN outperforms all other methods in a scenario composed of textural images and small-size datasets. Besides achieving higher accuracy, the proposed method has several advantages over the shallow and deep related methods. The LMFCN approach requires

Table 4: Average balanced accuracy and standard deviation for seven methods considering a multiclass scenario.

Dataset	Method	Training	Validation	Test
BACH	LMFCN	0.9690 ± 0.0348	0.7170 ± 0.0210	0.6854 ± 0.0278
	CNN	0.8390 ± 0.0184	0.6788 ± 0.0406	0.6444 ± 0.0187
	GLCM	0.4081 ± 0.0183	0.4003 ± 0.0753	0.3637 ± 0.0303
	LBP	0.4643 ± 0.0334	0.4871 ± 0.0729	0.3951 ± 0.0896
	PFTAS	0.6229 ± 0.0146	0.5690 ± 0.0464	0.6043 ± 0.0235
	ResNet18	1.0000 ± 0.0000	0.6056 ± 0.0578	0.5671 ± 0.0300
BreastHist	InceptionV3	0.4628 ± 0.0807	0.4269 ± 0.0904	0.4057 ± 0.0261
	LMFCN	0.9688 ± 0.0285	0.8125 ± 0.0321	0.7895 ± 0.0162
	CNN	0.8032 ± 0.0319	0.7347 ± 0.0238	0.7040 ± 0.0307
	GLCM	0.4697 ± 0.0111	0.4218 ± 0.0212	0.4091 ± 0.0050
	LBP	0.9235 ± 0.0084	0.5373 ± 0.0174	0.5218 ± 0.0147
	PFTAS	0.9537 ± 0.0086	0.6643 ± 0.0109	0.6768 ± 0.0147
CRC	ResNet18	1.0000 ± 0.0000	0.5933 ± 0.0341	0.5834 ± 0.0213
	InceptionV3	0.1848 ± 0.0182	0.1745 ± 0.0194	0.1767 ± 0.0123
	LMFCN	0.9834 ± 0.0127	0.9379 ± 0.0065	0.9338 ± 0.0073
	CNN	0.7209 ± 0.2614	0.3946 ± 0.0453	0.3901 ± 0.0402
	GLCM	0.6062 ± 0.0057	0.5960 ± 0.0069	0.6049 ± 0.0084
	LBP	0.6148 ± 0.0045	0.6086 ± 0.0270	0.6144 ± 0.0123
CRC	PFTAS	0.8359 ± 0.0058	0.8271 ± 0.0152	0.8297 ± 0.0097
	ResNet18	0.9506 ± 0.0033	0.5070 ± 0.0141	0.5066 ± 0.0138
	InceptionV3	0.1509 ± 0.0036	0.1546 ± 0.0075	0.1477 ± 0.0082

few data to train a sequence of CLs and a large-margin discriminant properly. In the experiments comparing the LMFCN with equivalent CNNs, using a 16-dimensional latent representation, the LMFCN achieved training stability and high accuracy within 20 training epochs. The CNNs needed 100 epochs to achieve comparable performance.

The latent representation learned by the LMFCN, constrained to a dimensionality similar to shallow methods (59- and 162-dimensional), is more discriminant than LBP, PFTAS, and GLCM. As a result, the LMFCN achieved higher balanced accuracy than the compared methods. Furthermore, the LMFCN obtained a competitive accuracy even with a 16-dimensional latent representation. Moreover, the improvement achieved by increasing the latent representation to 59 dimensions is meaningful, with a slight gain when increasing it up to 169.

The LMFCN reduces the number of SVs over the training epochs. Fewer SVs imply that the computational effort reduces, as the main term of the loss function depends on the number of SVs. Compared to equivalent CNNs, the LMFCN has an extra cost related to the computation of kernel \mathbf{K} , distance matrix \mathbf{D} , and the quadratic optimization problem solved by the SMO algorithm, which have computational complexities $O(n^2)$, $O(n^2)$, and $O(n^3)$, respectively. CNNs have a training complexity proportional to the number of instances and weights. They have more parameters, use more instances in the backpropagation, and take more epochs to converge than the LMFCN.

Data imbalance is burdensome for training machine learning algorithms. However, the LMFCN deals well with imbalanced classes in two-class and multi-class scenarios, achieving competitive performance even on highly imbalanced OVA subproblems. Furthermore, the computational effort is not extremely high, given the reduced number of epochs to train the model at each subproblem.

5. Conclusion

This paper proposed a large-margin representation learning approach that is made up of convolutional layers and a large-margin discriminant. As a result, the LMFCN achieved competitive accuracy compared to CNNs with similar architecture while reducing the computational cost. Furthermore, the LMFCN achieved training stability in a few epochs thanks to using

only the SVs in the backpropagation algorithm. These achievements were possible using a large-margin discriminant, which replaces the fully connected and softmax layers of conventional CNNs. As a result, an SVM with an RBF kernel can produce complex margins, avoiding expensive refinement of the latent representation. This way, the FCNs do not need to suffer drastic updates. Unlike the handcrafted features extractors, the LMFCN has more adaptability given its consistent results in different datasets. The most problematic scenario for the LMFCN is the multiclass classification task, but it showed promising results despite the overhead caused by the OVA approach. We used a few epochs per subproblem to alleviate the computational cost of using several models, keeping the cost of our approach similar to the CNNs. In conclusion, the LMFCN has advantages in computational cost and classification performance over the compared methods for small datasets of textural images.

Acknowledgments

This work was funded by the Natural Sciences and Engineering Research Council of Canada (NSERC) Grant RGPIN 2016-04855.

References

- [1] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *IEEE/CVF Conf Comp Vis Patt Recog*, 2016, pp. 770–778.
- [2] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the inception architecture for computer vision, in: *IEEE/CVF Conf Comp Vis Patt Recog*, 2016, pp. 2818–2826.
- [3] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei, ImageNet Large Scale Visual Recognition Challenge, *Intl Journal of Computer Vision* 115 (2015) 211–252. doi:10.1007/s11263-015-0816-y.
- [4] V. Andrearczyk, P. F. Whelan, Using filter banks in convolutional neural networks for texture classification, *Pattern Recognition Letters* 84 (2016) 63–69. doi:https://doi.org/10.1016/j.patrec.2016.08.016.
- [5] M. Cimpoi, S. Maji, A. Vedaldi, Deep filter banks for texture recognition and segmentation, in: *IEEE Conf Comp Vis Patt Recog*, 2015, pp. 3828–3836.
- [6] A. Humeau-Heurtier, Texture feature extraction methods: A survey, *IEEE Access* 7 (2019) 8975–9000.
- [7] S. Chopra, R. Hadsell, Y. LeCun, Learning a similarity metric discriminatively, with application to face verification, in: *IEEE CS Conf Comp Vis Patt Recog*, volume 1, 2005, pp. 539–546 vol. 1. doi:10.1109/CVPR.2005.202.
- [8] S. Kim, D. Kim, M. Cho, S. Kwak, Proxy anchor loss for deep metric learning, in: *IEEE/CVF Conf Comp Vis Patt Recog*, 2020.
- [9] X. Wang, Y. Hua, E. Kodirov, G. Hu, R. Garnier, N. M. Robertson, Ranked list loss for deep metric learning, in: *IEEE/CVF Conf Comp Vis Patt Recog (CVPR)*, 2019, pp. 5207–5216.
- [10] Y. Movshovitz-Attias, A. Toshev, T. K. Leung, S. Ioffe, S. Singh, No fuss distance metric learning using proxies, in: *IEEE Intl Conf on Comp Vis*, 2017, pp. 360–368.
- [11] N. Azire, S. Todorovic, Ensemble deep manifold similarity learning using hard proxies, in: *IEEE/CVF Conf Comp Vis Patt Recog*, 2019, pp. 7299–7307.
- [12] P. M. Roland Kwitt, Stex, salzburg texture image database (stex), <https://wavelab.at/sources/STex/>, 2019. Accessed: 2022-03-28.
- [13] F. A. Spanhol, L. S. Oliveira, C. Petitjean, L. Heutte, Breast cancer histopathological image classification using convolutional neural networks, in: *Intl Joint Conf Neural Networks*, 2016, pp. 2560–2567.
- [14] G. Aresta, T. Araújo, S. Kwok, et al., Bach: Grand challenge on breast cancer histology images, *Medical Image Analysis* 56 (2019) 122–139. doi:https://doi.org/10.1016/j.media.2019.05.010.
- [15] J. N. Kather, C.-A. Weis, F. Bianconi, S. M. Melchers, L. R. Schad, T. Gaiser, A. Marx, F. G. Zöllner, Multi-class texture analysis in colorectal cancer histology, *Scientific Reports* 6 (2016) 27988. doi:10.1038/srep27988.