

A NOVEL GRAPH-THEORETIC DEEP REPRESENTATION LEARNING METHOD FOR MULTI-LABEL REMOTE SENSING IMAGE RETRIEVAL

Gencer Sumbul and Begüm Demir

Faculty of Electrical Engineering and Computer Science, Technische Universität Berlin, Germany

ABSTRACT

This paper presents a novel graph-theoretic deep representation learning method in the framework of multi-label remote sensing (RS) image retrieval problems. The proposed method aims to extract and exploit multi-label co-occurrence relationships associated to each RS image in the archive. To this end, each training image is initially represented with a graph structure that provides region-based image representation combining both local information and the related spatial organization. Unlike the other graph-based methods, the proposed method contains a novel learning strategy to train a deep neural network for automatically predicting a graph structure of each RS image in the archive. This strategy employs a region representation learning loss function to characterize the image content based on its multi-label co-occurrence relationship. Experimental results show the effectiveness of the proposed method for retrieval problems in RS compared to state-of-the-art deep representation learning methods. The code of the proposed method is publicly available at <https://git.tu-berlin.de/rsim/GT-DRL-CBIR>.

Index Terms— Multi-label image retrieval, graph-theoretic representation learning, deep learning, remote sensing

1. INTRODUCTION

Multi-label content-based image retrieval (CBIR) methods aim to retrieve remote sensing (RS) images similar to a given query image by exploiting training images annotated by multi-labels. Development of effective CBIR methods has recently attracted great attention in RS. As an example, in [1] a sparse reconstruction-based multi-label RS image retrieval method that considers a measure of label likelihood is introduced. In [2], fully convolutional networks are introduced for multi-label RS images to extract descriptors of image regions in the content of CBIR. Recently, deep representation learning (DRL) methods based on a triplet loss function are found very popular for CBIR problems due to their intrinsic characteristic to model similarities of images. These methods employ image triplets (each of which includes anchor, positive and negative images), aiming to learn a metric space where the distance between the positive and the anchor images is minimized while that between the negative and anchor

images is maximized. In [3], triplet loss is employed with convolutional neural networks (CNN) to learn an embedding space for hash code generation of RS images. The use of triplet loss function requires an accurate selection of image triplets. A simple strategy is to define triplets from an existing training set of labeled images. However, such strategy does not guarantee the selection of the most informative images to the anchor, and thus can result in limited CBIR performance particularly when images annotated by multi-labels are available. In addition, the triplet selection based DLR methods do not take into account the co-occurrence relationships of land-cover classes present in an RS image. However, modeling these relationships is crucial for an accurate CBIR. This problem can be addressed by using graphs, which capture both region characteristics and the spatial relationships among the regions. In [4], a semi-supervised graph-theoretic method is introduced to model inherent correlation of multi-labels by a correlated label propagation algorithm. The performance of this approach depends on the hand-crafted features to represent each image region. Recently, in [5] region graph-based image representations are utilized to model the similarity of image pairs via a siamese graph CNN in the context of DRL. This method learns a metric space based on only pairwise image similarities, which may not be sufficient to model the complex information content of RS images for CBIR problems.

To address the above-mentioned issues, in this paper we propose a graph-theoretic deep representation learning method that does not require image pairs and triplets. The proposed method models multi-label co-occurrence relationships based on a novel region representation learning loss function.

2. THE PROPOSED GRAPH-THEORETIC DEEP REPRESENTATION LEARNING METHOD

Let $\mathcal{X} = \{x_1, \dots, x_I\}$ be an archive that includes I images, where x_j is the j^{th} RS image in the archive \mathcal{X} . We assume that a training set $\mathcal{T} \subset \mathcal{X}$ that consists of labeled images is available. Each image in \mathcal{T} is associated to pixel-based labels from a label set $\mathcal{B} = \{l_1, \dots, l_C\}$. Let m_i be the land-cover map of the image $x_i \in \mathcal{T}$ (m_i and x_i have the same pixel sizes and thus each pixel in m_i represents the label of the

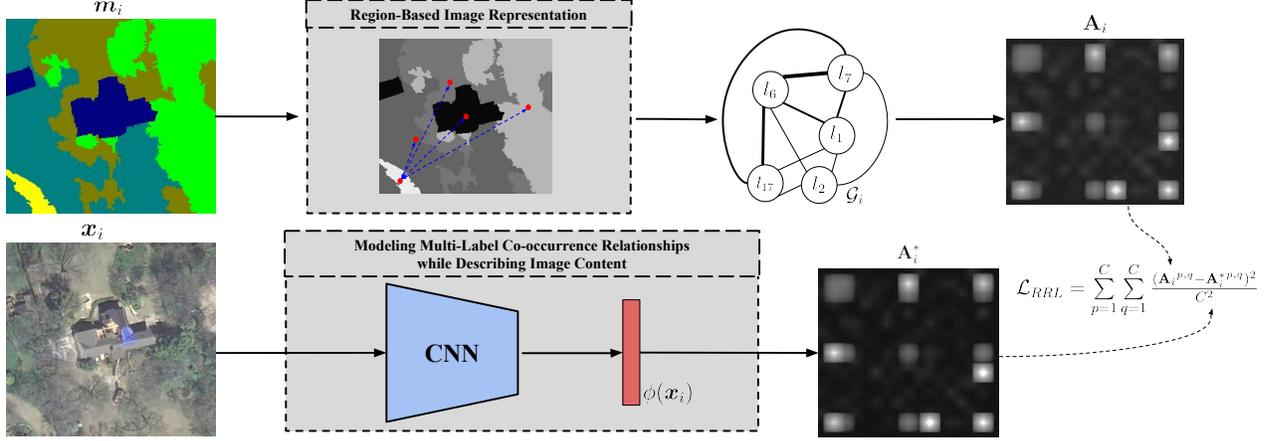


Fig. 1. Illustration of the proposed graph-theoretic deep representation learning method.

corresponding pixel in x_i). The set of all labels associated to x_i are defined by a binary vector $\mathbf{y}_i \in \{0, 1\}^C$, where each element of \mathbf{y}_i indicates the presence or absence of label $l_n \in \mathcal{B}$.

The proposed method aims to model co-occurrence relationship of multiple classes present in each image in the archive. To this end, each training image x_i is represented with a graph structure, which provides region-based image representation (where each region is associated with a land-cover class). The proposed method includes a novel learning strategy to automatically predict the corresponding graph structure of any image in the archive, while describing the complex content of each image. To this end, we exploit a convolutional neural network (CNN). However, the proposed learning strategy can be injected to any deep neural network. Fig. 1 shows a general overview of the proposed method, which is explained in detail in the following.

During training, to describe regions associated to classes present in each training image, the proposed method first constructs a graph structure, where the nodes represent the image region properties and the edges represent the spatial relationship among the regions. Let $\mathcal{G}_i = (E_i, V_i, \mathbf{W}_i)$ be the graph associated to the image x_i . E_i is the set of graph edges, V_i is the set of nodes and $\mathbf{W}_i \in \mathbb{R}^{C \times C}$ is the weight matrix of the graph. Each node represents a region associated to a class of the image (i.e., $l_n \in \mathbf{y}_i$). The weight $\mathbf{W}_i^{p,q}$ between l_p and l_q $\mathbf{W}_i^{p,q} = 1$ if $l_p \in \mathbf{y}_i, l_q \in \mathbf{y}_i$, and otherwise $\mathbf{W}_i^{p,q} = 0$.

By this way, all the class relationships of x_i are modeled with same importance. However, class relationships can be subject to different levels of importance based on the characteristics of each region and its spatial relationship with the other regions. As an example, the relationships between an image region and its neighbors are more important than those between non-neighbor regions. In detail, if two neighbor regions cover most of the image content, their relationship plays the most significant role for accurately modeling the multi-label co-occurrence relationship. To address this issue, we

define a weight for the edge $\mathbf{W}_i^{p,q}$ as follows:

$$\mathbf{W}_i^{p,q} = \frac{s(l_p; \mathbf{m}_i) \times s(l_q; \mathbf{m}_i)}{N_s} \times \left(1 - \frac{d(l_p, l_q; \mathbf{m}_i)}{N_d}\right) \quad (1)$$

where $s : \mathcal{B} \mapsto \mathbb{N}$ is a function that maps a class label into the size of the region associated to the class, $d : \mathcal{B} \times \mathcal{B} \mapsto \mathbb{N}$ is a function that maps the pairs of class labels into the distance between the centers of their regions associated to the corresponding classes. N_s and N_d are the maximum values of the functions s and d , respectively. By this way, if the regions are close to each other and their sizes are large, the weights assigned to the corresponding edges in the graph \mathcal{G}_i will be high. After obtaining a graph for each training image, the characteristics and spatial arrangements of image regions are represented with an adjacency matrix $\mathbf{A}_i \in \mathbb{R}^{C \times C}$ where $\mathbf{A}_i^{p,q} = \mathbf{W}_i^{p,q}$ if an edge exists between the nodes V_i^p and V_i^q in the graph \mathcal{G}_i , $\mathbf{A}_i^{p,q} = 0$ otherwise.

To model multi-label co-occurrence relationship of any image in the archive, the proposed learning strategy consists of region-based image representation learning and image characterization. Let $\phi : \theta, \mathcal{X} \mapsto \mathbb{R}^\gamma$ be any type of CNN that maps the image x_i to γ -dimensional image descriptor, where θ is the set of CNN parameters. The region-based image representation learning is achieved by the prediction of the adjacency matrix based on the image descriptor. To this end, the characterization of x_i is performed based on the considered CNN to model the multi-label co-occurrence relationship of x_i in the adjacency matrix \mathbf{A}_i . The prediction of the adjacency matrix is achieved by a fully connected layer that takes the image descriptor $\phi(x_i)$ and produces the vectorized form of the reconstructed adjacency matrix. To train the proposed method, we define a novel region representation learning loss $\mathcal{L}_{\mathcal{RRL}}$ function as follows:

$$\mathcal{L}_{\mathcal{RRL}} = \sum_{x_i \in \mathcal{T}} \sum_{p=1}^C \sum_{l=1}^C \frac{(\mathbf{A}_i^{p,q} - \mathbf{A}_i^{*p,q})^2}{C^2}. \quad (2)$$

The proposed loss function allows to describe the content of

Table 1. Mean average precision (mAP) obtained for the DLRSD and BigEarthNet-S2 archives.

Method	Benchmark Archive	
	DLRSD	BigEarthNet-S2
SNN (random) [11]	66.5%	83.9%
SNN (batch-all) [11]	68.0%	88.6%
SNN (hard) [11]	70.7%	88.3%
SGCN [5]	70.1%	87.8%
Proposed Method	84.3%	92.1%

an RS image based on the multi-label co-occurrence information to achieve the region-based image representation learning. After an end-to-end training of the whole neural network by minimizing the region representation learning loss and thus learning the network parameters $\theta^* = \arg \min_{\theta} \mathcal{L}_{\mathcal{R}\mathcal{R}\mathcal{L}}$, the proposed method extracts the descriptors $\{\phi(\mathbf{x}_j; \theta^*)\}$ of the images in the archive \mathcal{X} . To perform CBIR, the proposed method retrieves RS images from the archive similar to a given query image \mathbf{x}_q by comparing $\phi(\mathbf{x}_q)$ with each element of the set $\{\phi(\mathbf{x}_j; \theta^*)\}$.

It is worth noting that the proposed method considers \mathbf{m}_i of $\mathbf{x}_i \in \mathcal{T}$ (i.e., pixel-level labels of training images) is available for the training phase. In the case that training images are annotated by image-level multi-labels instead of pixel-level labels, \mathbf{m}_i can be obtained by using a weakly-supervised semantic segmentation that exploits only image-level annotations as explained in [6].

3. EXPERIMENTAL RESULTS

Experiments were conducted on the DLRSD [7] and the BigEarthNet-S2 [8] benchmark archives. The DLRSD archive is the extension of the UC Merced archive [9] that includes 2,100 aerial images, each of which has the size of 256×256 pixels with a spatial resolution of 30 cm. The DLRSD archive also includes pixel labels defined in [4]. To perform experiments, we split the DLRSD archive into training (80%) and test (20%) sets. The large-scale BigEarthNet-S2 benchmark archive consists of 590,326 Sentinel-2 images. Each image in BigEarthNet-S2 has been annotated with multi-labels from the 2018 CORINE Land Cover (CLC) database. In this paper, we first extracted the CLC land cover map of each image and then exploited it based on the 19 classes nomenclature presented in [10]. To perform experiments, we first selected the 74,716 BigEarthNet-S2 images acquired over Serbia and then divided them into training (52%), validation (24%) and test (24%) sets. To select query images, the training set of the DLRSD archive and the validation set of the BigEarthNet-S2 archive were used, while images were retrieved from the test set for both archives. In the experiments, we exploited the DenseNet model [12] at the depth of 121. We trained our method for 100 epochs by using the Adam optimizer. We compared our method with

siamese neural networks (SNNs) trained with triplet loss [11] and siamese graph convolution network (SGCN) trained with contrastive loss [5]. For SNN, we utilized random, batch-all and hard sampling techniques in the experiments. The results are denoted as SNN (random), SNN (batch-all) and SNN (hard). The reader is referred to [13] for the details of these techniques. The same training procedure and the same backbone with the proposed method were used for all experiments. For SGCN, we employed the same graph formation and parameter values given in [5]. To obtain CBIR results, chi-square distance is utilized to compare image descriptors. Experimental results are provided in terms of normalized discounted cumulative gains (NDCG), mean average precision (mAP) and average cumulative gains (ACG) [14]. Table 1 shows the mAP results obtained on both archives. By assessing the table, one can observe that the proposed method leads to the highest mAP scores compared to the SNN with all types of sampling techniques and SGCN. As an example, the proposed method provides almost 18% higher and more than 8% higher mAP scores for DLRSD and BigEarthNet-S2 archives, respectively, compared to the SNN (random). This shows that modeling multi-label co-occurrence of an RS image by the proposed method improves the CBIR performance compared to the SNN, in which multi-label dependencies present in an image have not been considered. As an other example, the proposed method provides almost 14% higher mAP score for the DLRSD archive compared to the SNN (hard). In addition, the proposed method leads to more than 14% higher and almost 5% higher mAP scores for DLRSD and BigEarthNet-S2 archives, respectively, compared to the SGCN, which is one of the state-of-the-art graph-based DRL methods for CBIR. These results show that, without a need for pair or triplet selection, the proposed method characterizes the image similarity much more accurately compared to the triplet and contrastive loss based DRL methods. Fig. 2 shows the mAP, ACG and NDCG results for the DLRSD archive under different numbers of retrieved images. By analyzing the figure, one can see that increasing the number of retrieved images does not change our conclusion. As an example, the proposed method outperforms SCNN and SGCN by almost 20% in NDCG for the DLRSD archive when the number of retrieved images is 100. It is worth noting that the promising CBIR performance of our method relies on: i) accurately predicting the graph structures of images in the archive; and ii) defining image descriptors based on the co-occurrence relationship of land-cover classes present in each image.

4. CONCLUSION

In this paper, we have presented a novel graph-theoretic deep representation learning method for multi-label RS image retrieval problems. The effectiveness of our method relies on the efficient use of a novel learning strategy that contains a region representation learning loss function (which allows

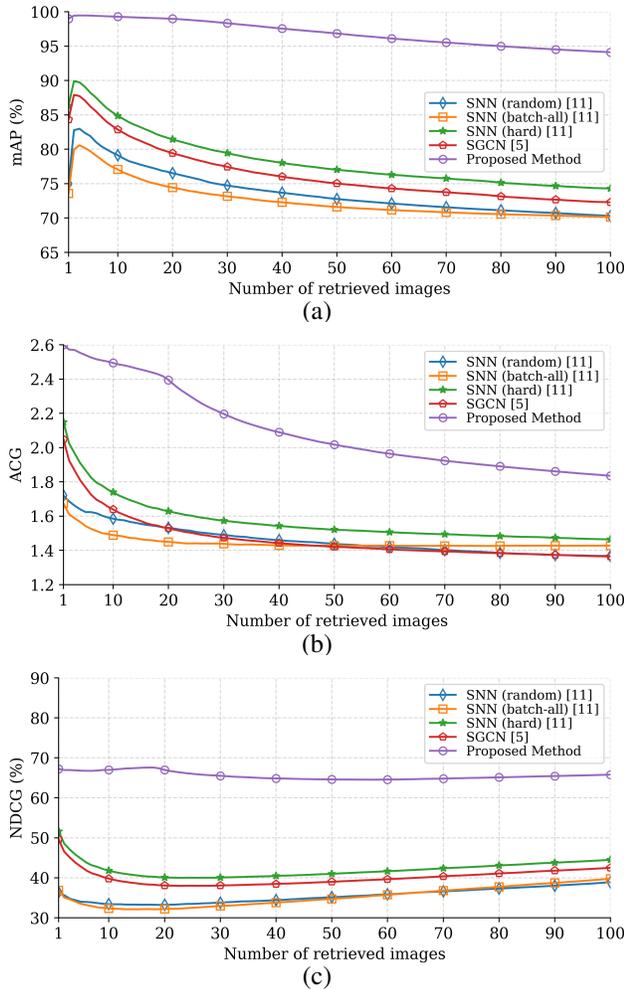


Fig. 2. (a) Mean average precision (mAP); (b) Average cumulative gains (ACG) and (c) Normalized discounted cumulative gains (NDCG) versus the number of retrieved images obtained for the DLRSD archive.

to model an RS image content based on the multi-label co-occurrence relationships). Experimental results show the success of the proposed method compared to state-of-art deep representation learning methods.

It is worth noting that the proposed method requires training images annotated by pixel-level labels. Such class labels can be attained through publicly available thematic products. However, class labels available through the thematic products can be noisy (incomplete, outdated, etc.), and thus their direct use may result in an uncertainty in the DL models and thus uncertainty in the CBIR performance. As a future work, we plan to develop label-noise robust graph-theoretic deep representation learning methods.

5. ACKNOWLEDGEMENTS

This work was supported by the European Research Council under the ERC Starting Grant BigEarth-759764.

6. REFERENCES

- [1] O. E. Dai, B. Demir, B. Sankur, and L. Bruzzone, "A novel system for content-based retrieval of single and multi-label high-dimensional remote sensing images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 11, no. 7, pp. 2473–2490, 2018.
- [2] Z. Shao, W. Zhou, X. Deng, M. Zhang, and Q. Cheng, "Multilabel remote sensing image retrieval based on fully convolutional network," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 13, pp. 318–328, 2020.
- [3] S. Roy, E. Sangineto, B. Demir, and N. Sebe, "Metric-learning-based deep hashing network for content-based retrieval of remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 18, no. 2, pp. 226–230, 2021.
- [4] B. Chaudhuri, B. Demir, S. Chaudhuri, and L. Bruzzone, "Multilabel remote sensing image retrieval using a semisupervised graph-theoretic method," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 2, pp. 1144–1158, 2018.
- [5] U. Chaudhuri, B. Banerjee, and A. Bhattacharya, "Siamese graph convolutional network for content based remote sensing image retrieval," *Comput. Vis. Image Understand.*, vol. 184, pp. 22–30, 2019.
- [6] L. Chan, M. S. Hosseini, and K. N. Plataniotis, "A comprehensive analysis of weakly-supervised semantic segmentation in different image domains," *Int. J. Comput. Vis.*, 2020.
- [7] Z. Shao, K. Yang, and W. Zhou, "Performance evaluation of single-label and multi-label remote sensing image retrieval using a dense labeling dataset," *Remote Sens.*, vol. 10, no. 6:964, 2018.
- [8] G. Sumbul, M. Charfuelan, B. Demir, and V. Markl, "BigEarthNet: A large-scale benchmark archive for remote sensing image understanding," *IEEE Intl. Geosci. Remote Sens. Symp.*, pp. 5901–5904, 2019.
- [9] Y. Yang and S. Newsam, "Geographic image retrieval using local invariant features," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 2, pp. 818–832, 2013.
- [10] G. Sumbul, A. d. Wall, T. Kreuziger, F. Marcelino, H. Costa, P. Benevides, M. Caetano, B. Demir, and V. Markl, "BigEarthNet-MM: A large scale multi-modal multi-label benchmark archive for remote sensing image classification and retrieval," *arXiv preprint arXiv:2105.07921*, 2021.
- [11] Florian S., Dmitry K., and James P., "Facenet: A unified embedding for face recognition and clustering," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2015, pp. 815–823.
- [12] G. Huang, Z. Liu, L. v. d. Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017, pp. 2261–2269.
- [13] A. Hermans, L. Beyer, and B. Leibe, "In Defense of the Triplet Loss for Person Re-Identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [14] Z. Zhang, Q. Zou, Y. Lin, L. Chen, and S. Wang, "Improved deep hashing with soft pairwise similarity for multi-label image retrieval," *IEEE Trans. Multimedia*, vol. 22, no. 2, pp. 540–553, 2020.