The Institution of Engineering and Technology WILEY

ORIGINAL RESEARCH

OSAP-Loss: Efficient optimization of average precision via involving samples after positive ones towards remote sensing image retrieval



¹School of Computer Science and Technology, Wuhan University of Science and Technology, Wuhan, China

²Hubei Province Key Laboratory of Intelligent Information Processing and Real-time Industrial System, Wuhan University of Science and Technology, Wuhan, China

³Nanyang Technological University, Singapore, Singapore

⁴School of Computer Science, Wuhan University, Wuhan, China

⁵Department of Electrical Engineering, Industrial Technology Research Institute, National Tsing Hua University, Hsinchu, China

Correspondence

Xin Xu, School of Computer Science and Technology, Wuhan University of Science and Technology, Wuhan, China. Email: xuxin@wust.edu.cn

Funding information

National Natural Science Foundation of China, Grant/Award Numbers: U1803262, 62176191, 62171325; Nature Science Foundation of Hubei Province, Grant/Award Number: 2022CFB018; Hubei Province Key Laboratory of Intelligent Information Processing and Real-time Industrial System, Grant/Award Number: ZNXX2022001

Abstract

In existing remote sensing image retrieval (RSIR) datasets, the number of images among different classes varies dramatically, which leads to a severe class imbalance problem. Some studies propose to train the model with the ranking-based metric (e.g., average precision [AP]), because AP is robust to class imbalance. However, current AP-based methods overlook an important issue: only optimising samples ranking before each positive sample, which is limited by the definition of AP and is prone to local optimum. To achieve global optimisation of AP, a novel method, namely Optimising Samples after positive ones & AP loss (OSAP-Loss) is proposed in this study. Specifically, a novel superior ranking function is designed to make the AP loss differentiable while providing a tighter upper bound. Then, a novel loss called Optimising Samples after Positive ones (OSP) loss is proposed to involve all positive and negative samples ranking after each positive one and to provide a more flexible optimisation strategy for each sample. Finally, a graphics processing unit memory-free mechanism is developed to thoroughly address the non-decomposability of AP optimisation. Extensive experimental results on RSIR as well as conventional image retrieval datasets show the superiority and competitive performance of OSAP-Loss compared to the state-of-the-art.

KEYWORDS

computer vision, image retrieval, metric learning

1 | INTRODUCTION

In recent years, with the investment in scientific research of satellite remote sensing technology, a large amount of relevant remote sensing data has been produced. As a result, numerous application research studies have been derived, among which remote sensing image processing applications have attracted wide attention. At the same time, these applications are in urgent need of effective remote sensing image retrieval (RSIR) techniques, which refers to the task of ranking semantically matched or similar images in a large remote sensing image database based on their relevance to the query [2–4]. RSIR has

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

^{© 2023} The Authors. CAAI Transactions on Intelligence Technology published by John Wiley & Sons Ltd on behalf of The Institution of Engineering and Technology and Chongqing University of Technology.

been studied for many years, and its core is how to quickly and accurately find the target image from a large number of remote sensing images. The performance of the predicted results in RSIR is often evaluated using ranking-based metrics, for example, the standard recall at k (Recall@k) [5], the normalised discounted cumulative gain (NDGG) [6], average precision (AP) [7, 8], and spearman coefficient [9]. Nevertheless, it is challenging to design a high-performance method for RSIR due to the complex and variable geographic information contained in remote sensing images.

With the continuous research and development of deep learning, end-to-end training in RSIR has become the de facto choice. Although deep neural networks (DNNs) combined with deep metric learning have achieved some success for RSIR, they are tenuous to class imbalance [10] (i.e. the number of images among different classes varies dramatically in RSIR datasets), which will lead to performance degradation. As the ranking evaluation metrics are robust to class imbalance according to their ranking-based error definition [11]. It has been proposed to directly optimise the well-defined ranking-based evaluation metrics to deal with the class imbalance problem. This intuition appears attractive, but it is notoriously difficult due to two major challenges posed by these evaluation metrics: (i) non-differentiability and (ii) non-decomposability. To tackle these problems, intuitively surrogate loss functions such as 0-1 loss [12, 13], the area under the ROC curve [14, 15], and cross entropy [16] were utilised extensively. These loss functions are decomposable, that is, they can be decomposed over every training sample. However, these attempts were merely made by optimising a structured hinge-loss upper bound to such evaluation metrics [6, 8] rather than the metrics themselves. Recently, there have been many studies using asymptotic methods to directly optimise, for example, binning approaches [17-21], neural networks [22], and blackbox methods [23, 24].

These methods provide an elegant upper bound for such metrics, but they are generally coarse approximations. Consequently, some efforts have attempted to design more accurate and smooth approximations [17–22, 25]. All studies conducted over many years have greatly promoted the optimisation of ranking-based evaluation metrics for vision tasks, such as conventional image retrieval (CIR) [26–30], face recognition [31–33], person/vehicle re-identification [34–39], and object detection [11] etc. Among these evaluation metrics, the AP is a pivotal evaluation metric in RSIR.

Therefore, we consider AP optimisation as a solution to overcome the class imbalance problem in RSIR. However, there is an important issue with AP optimisation that has been overlooked by previous research: only optimising samples ranking before each positive sample, resulting in that all sample information after each positive sample is not fully exploited. This issue easily makes AP optimisation fall into a local optimum. Moreover, this issue has rarely received attention so far, and it is the first to focus on this issue in this paper. In addition, current AP-based methods adopt one special gradient strategy, which may assign a smaller gradient to positive or negative samples with higher scores referenced to a positive sample with a lower score. Therefore, it should have different gradients for different positive and negative samples, which facilitates adaptive optimisation and improves the robustness of the trained model. As shown in Figure 1, the intra-class diversity and inter-class similarity of remote sensing images are exceptionally significant, which poses a great challenge to the robustness of the learnt model. These issues and challenges affect the retrieval results of AP optimisation for RSIR.

In this paper, we tackle the above issues of directly optimising AP with the stochastic gradient descent (SGD) methods for RSIR, which not only address the non-differentiable and non-decomposable nature of ranking to AP, but also include the



FIGURE 1 Sample images of four different categories (e.g. baseball field, basketball court, freeway, and runway) from PatternNet [1] dataset. (a) Baseball field and (b) basketball court images, where the images with the same category show significant diversity. And (c) freeway and (d) runway images, where extreme similarity exists between the images of different categories.

specific optimisation limited by the AP definition. Existing APbased methods have investigated and achieved some success in the non-differentiability and non-decomposability of AP, but improvement is still limited (detailed theoretical analysis is shown in Section 4). Therefore, a more accurate and effective optimisation of AP loss remains an open and explorable issue.

To this end, we propose a novel method, namely Optimising Samples after positive ones & AP loss (OSAP-Loss) for RSIR, which provides an efficient training objective for the optimisation and improvement of AP. Specifically, OSAP-Loss includes three components: \mathcal{L}_{AP}^{SRF} , \mathcal{L}_{OSP} , and memory-free mechanism (MFM). First, to address the non-differentiability of AP and avoid the vanishing gradient of some samples, a novel superior ranking function (SRF) is proposed to enable AP loss differentiable while providing a tighter upper bound, thereby ensuring that it can be directly optimised using standard gradient descent methods, that is, \mathcal{L}_{AP}^{SRF} . Then, we introduce a novel loss, called Optimising Samples after Positive ones (OSP) loss (i.e. \mathcal{L}_{OSP}), to reduce intra-class variability and inter-class similarity for further optimising \mathcal{L}_{AP}^{SRF} . More importantly, the \mathcal{L}_{OSP} considers the ranking of all positive and negative samples after each positive one and offers a more flexible optimisation approach for each sample, which maintains a margin between positive and negative samples. In this case, the initial version of OSAP-Loss (i.e. \mathcal{L}_{OSAP}) is formulated, that is, by linearly combining \mathcal{L}_{AP}^{SRF} and \mathcal{L}_{OSP} . To achieve a small deviation in AP optimisation, a graphics processing unit (GPU) MFM is developed to sidestep the hardware constraints of the GPU memory, which could thoroughly address the non-decomposability in AP optimisation. The MFM is inspired by the instance-level retrieval work [20], and we are the first to apply this technology to the field of RSIR. It is possible to use a very large batch size with several thousand high-resolution images on a single GPU. At last, \mathcal{L}_{OSAP}^{MFM} is derived from \mathcal{L}_{OSAP} with MFM. This straightforward yet effective loss function can improve the robustness of the trained model for large-scale RSIR tasks. Meanwhile, we provide a theoretical analysis and discussion to verify the superiority and effectiveness of our proposed method.

In summary, the major contributions of this paper are as follows:

- We propose shifting ranking orders rather than modifying the metric learning loss and minimising distances in the embedding space. In this paper, we propose a new method, named OSAP-Loss, which provides an efficient training objective for the optimisation and improvement of AP. To avoid the vanishing gradient of some samples, we design a novel SRF that is differentiable for the AP loss (i.e. \$\mathcal{L}_{AP}^{SRF}\$) while providing a tighter upper bound, thus ensuring that it can be directly optimised using standard gradient descent methods.
- 2) To reduce intra-class variability and inter-class similarity for further optimising \mathcal{L}_{AP}^{SRF} , we design a novel loss, for example, OSP loss, to involve all samples after positive ones while offering a more flexible optimisation approach for each sample for global ranking optimisation. In

addition, we are the first to develop a GPU MFM to thoroughly address the non-decomposability of AP by sidestepping memory constraints of GPU in RSIR. The MEM enables the training of arbitrarily sized batches, which guarantees the global optimisation of AP and thereby improves the overall ranking performance.

3) A theoretical analysis is provided to show the superiority and effectiveness of our method. More importantly, we conduct extensive experiments on six image retrieval benchmarks, that is, three RSIR benchmarks containing UCMD [40], NWPU-RESISC45 [41], and PatternNet [1], as well as three CIR benchmarks, including CUB-200-2011 [42], Stanford Online Product (SOP) [43], and INaturalist-2018 [44]. Experimental results show that our proposed OSAP-Loss is on par with or superior to the state-of-theart, which also demonstrates its effectiveness.

The rest of this paper is organised as follows. The relevant works on CIR, RSIR, deep metric learning and direct optimisation for the AP metric are discussed in Section 2. We introduce some preliminaries used in this work in Section 3 and describe the proposed method in Section 4. In Section 5, we elaborate experimental settings and extensive experimental results. Finally, conclusions are drawn in Section 6.

2 | RELATED WORK

In this section, we first review the related works of CIR in Section 2.1, and then we revisit the studies on RSIR in Section 2.2. Next, Section 2.3 provides the description and analysis of deep metric learning. Finally, we introduce the methods of direct optimisation for AP in deep metric learning in Section 2.4.

2.1 | Conventional image retrieval

CIR has been a fundamental and hot research topic in the field of information retrieval. It aims to find all images with relevant content to the query in the database. Existing efforts on CIR mainly focus on two categories: (1) the speed of retrieval and (2) the accuracy of retrieval. For the speed of retrieval, with the explosive growth of image content on the Internet, how to conduct fast and effective retrieval has emerged as major attention, where deep hashing [45-47] has become a leading technique for fast image retrieval. Previous studies mainly utilised hand-crafted image features like scale-invariant feature transform (SIFT) to learn hash functions for modelling data structures to preserve image similarities. Existing hashing methods directly extract features to learn hash codes or hash functions [48, 49] with the help of capable DNNs. For the accuracy of retrieval, it focuses on designing or learning more efficient image representations to achieve higher accuracy image retrieval performance. The earlier works aimed at obtaining compact image descriptors composed of multiple local features,

such as Fisher vectors [50] and vector of locally aggregated descriptors (VLAD) [51, 52]. In recent years, DNNs have made wonderful progress in CIR due to their powerful non-linear fitting and feature capturing capabilities. On the one hand, discriminative descriptors are generated by designing specific neural network structures. On the other hand, designing training objectives [43, 53–58] to achieve robust and generalised feature distributions is another highly typical research area, namely deep metric learning. It is worth noting that CIR can be extensively viewed as a problem of learning to rank. Meanwhile, the key to deep metric learning is to design a loss function that optimises a good ranking as opposed to classification, which is consistent with the objective of image retrieval. Therefore, image retrieval can be solved by combining deep metric learning with an appropriate ranking loss.

2.2 | Remote sensing image retrieval

RSIR is closely similar to CIR, except that the images retrieved are remote sensing images. In the last decade, RSIR has primarily focussed on the effective and discriminative feature extraction [2, 4, 59-62] and the dataset construction [1, 40, 41] of remote sensing image scenes. The field of feature extraction can be roughly divided into two categories: (1) hand-crafted features (e.g. SIFT [61], colour features [62], texture features [63], shape features [64], and spatial relationships [65] etc.); (2) deep features extracted by DNNs [2, 4, 59, 60]. Compared with hand-crafted features, it has become a consensus that deep learning can improve the accuracy of RSIR tasks, thanks to the tremendous feature capabilities of DNNs. Meanwhile, with the construction of research datasets in the field of remote sensing (e.g. RSIR datasets containing NWPU-RESISC45 [40], NWPUD [41], and PatternNet [1] etc.), which provides a large amount of training data for deep learning methods. Currently, the DNN-based RSIR method is to extract the features of remote sensing images with DNNs, and then train the model using classification loss [4, 66]. And deep metric learning methods have great potential in RSIR. The loss function plays a vital role in improving the discrimination and distribution of the learned features by DNNs. At present, some loss functions have been applied to RSIR with good results, including contrastive loss [67], triplet loss [68], N-Pair loss [69], Proxy NCA loss [70], lifted structured loss [43], and distribution structure learning loss (DSLL) [71]. However, these methods ignore the distribution of intra-class and the differences of inter-class. To this problem, Fan et al. [2] proposed distribution consistency loss (DCL) [2] to select multiple positive and negative samples from different classes to maintain intraclass compactness as well as inter-class differences by separating the negative samples from different classes by different distances. Despite these deep metric learning methods being novel and effective for RSIR, these loss functions belong to local optimisation and are challenging for high-resolution images in RSIR. As stated previously, image retrieval itself is a ranking problem, and the current loss function optimises merely the upper bound of the ranking loss. Therefore, a globally optimised loss function is necessary for RSIR in deep metric learning.

24682222, 0, Downloaded from https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cit2.12151, Wiley Online Library on [2505/2023]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

2.3 | Deep metric learning

Deep metric learning plays a critical role in many tasks, which has been studied for decades. It concentrates on constructing an effective feature space that effectively reflects the similarity or dissimilarity among images. The objective of deep metric learning is to minimise intra-class similarity while simultaneously maximising inter-class similarity. Metric learning has been extensively investigated by learning Mahalanobis distance functions [72] or projection matrix [73] before deep learning received widespread attention. Existing studies based on deep metric learning have mainly focussed on the designing of suitable loss functions to guide feature representation learning. For example, widely used pair-wise (e.g. contrastive loss [67]) and tuple-wise (e.g. triplet based loss [58, 68, 74-76] and n-tuplebased loss [69, 70, 77]) have been investigated for many years in the field of image retrieval. The essence of these losses is to construct the similarity structure between positive and negative samples by enlarging the inter-class distance and reducing the intra-class distance (e.g. Euclidean distance or Cosine distance). Nevertheless, these methods typically perform a local optimisation due to the fact that they act on a limited and fixed number of samples during the training phase. In addition, these methods need a series of iterative operations during training, including the repeatedly random sampling of challenging image pairs or tuples, especially time-consuming hard sample mining, as well as the calculation of losses and back-propagating gradients.

Moreover, a number of current studies [78-80] have pointed out that above-mentioned methods have some potential limitations resulting in the local optimisation. Firstly, it can be noticed that in most of these methods, only a proportion of informative samples is considered for constructing a similarity structure. As a result, numerous non-trivial samples are ignored, biasing the construction of the similarity structure. Secondly, intra-class data distribution properties are often ignored by most of these methods. Especially, these methods [43, 69, 75] attempt to pull the samples from the same class as closely as possible, and in extreme cases, compress the samples from the same class into a single point, which may easily degrade their similarity structure. Thirdly, these methods treat each sample with a specific gradient optimisation strategy, which forces the optimisation of the model to spend more capacity on some samples that have a minor effect on the similarity structure. Importantly, another often overlooked yet essential issue so far is that these loss functions only optimise the upper bound surrogate loss of true ranking loss [8, 23, 81-83]. Therefore, when using the ranking metric for evaluation, the minimum of the loss does not always correspond to the minimum of the true ranking loss (i.e. sub-optimal), namely the minimum of the loss does not guarantee that it corresponds to the minimum of the true ranking loss when using the ranking metric for evaluation. Furthermore, the common property of existing deep metric learning methods is that most of them are driven by minimising distance (maximising similarity) and thus overlook the importance of shifting ranking orders, which is crucial when evaluated with rank-based metrics. In this case, there have been extensive

research efforts [17, 20, 24, 83-85] focussing on the direct optimisation of ranking metric, such as AP.

2.4 Direct optimisation for AP

AP is a widely used evaluation metric in many tasks, such as object detection, image retrieval, person/vehicle reidentification and etc. Consequently, AP is a good, direct, and effective choice as an optimisation goal. Recently, some methods [17, 20, 24, 83–85] have been developed to directly optimise AP as the training objective. Nevertheless, there are two major challenges for optimising the approximation of AP: non-differentiability and non-decomposability. To address these two challenges, there are many studies in image retrieval proposed and dedicated to AP optimisation.

For the first challenge, it has been actually studied for many years. On the one hand, one solution is to find the smooth surrogate losses for AP approximation. The widely used surrogate losses for image retrieval are often based on the familiar losses (e.g. contrastive loss [67], triplet loss [68], quadruplet loss [86], or n-tuple loss [69, 70, 77]) to enforce local ranking. These methods only provide a very coarse upper bound for optimising the AP, and require a complex sampling strategy and some experience tricks to be effective. To this end, studying and designing a smooth upper bound on AP has received many progressive investigations [8, 23, 81-83]. Initially, some works have adapted structured SVM models [8, 87] to reduce the complexity of the corresponding loss-enhanced inference or to adjust to weak supervision [81]. After that, there are many AP optimisation methods, for example, using a large LSTM [22] or blackbox optimisation [24] to approximate the ranking step. On the other hand, another option is to design smooth approximations of the ranking function. The most typical approach is the soft-binning techniques [17-21] via smoothing discrete approximation of similarity scores. In recent years, another approach has been explored directly to approximate the nondifferentiable part of the ranking function, such as neural networks [22] or sigmoid functions [25, 84]. These approaches provide a more accurate approximation of AP by designing tight and smooth approximations of the ranking function. Although these two categories of methods provide applicable and smooth AP approximation, there are still some limitations. Specifically, the former methods provide subtle upper bounds, but it is mostly a coarse AP approximation; although the latter methods provide accurate AP approximation, they cannot guarantee to provide an upper bound for AP loss.

For the second challenge, the main reason lies in that the AP metric cannot be linearly decomposed into individual samples. Specifically, due to the limited available computational space, existing training methods are based on linearly dividing the whole samples into multiple batches, in which there will be a gap between the calculated AP of multiple batches and the AP of the whole samples. Consequently, it delivers an inconsistent AP gradient estimator. To address this challenge, effective batch sampling [5, 74, 88] or selecting informative samples [69, 88–90] makes the data distribution in a mini-batch

as close as possible to the whole data distribution. And, Wang et al. [91] proposed the cross-batch memory method to store learned features and compute the global objective based on the assumption of the slow drift in learned features. Nondecomposability has been effectively handled in AP optimisation by increasing the batch size with a brute-force effort [17, 20, 23, 25]. Although this way is effective in mitigating nondecomposability and improving model performance, it introduces an important overhead in the computation and memory of GPU, including two main steps for the computation of AP loss and the update of back-propagation gradients. Therefore, Ramzi et al. [85] tried to enable good performance for AP optimisation with a simple but effective loss \mathcal{L}_{calibr} , which uses small batches without introducing any overhead. Although these methods have been proposed to alleviate the decomposability gap and achieve significant success, it does not completely and essentially solve this issue.

Although the above works have been proposed to tackle these two major challenges in AP optimisation, the improvement of non-differentiability and non-decomposability is still limited. In addition, there is an overlooked issue of AP optimisation: only optimising samples ranking before each positive sample, which results in that all sample information after each positive sample is not fully exploited. In other words, existing AP based losses do not achieve genuine global ranking optimisation. In order to overcome these limitations, this paper proposes a new optimised loss function, named OSAP-Loss, to achieve global ranking optimisation for RSIR. Furthermore, we also discuss and analyse the superiority and effectiveness of our proposed method theoretically.

3 | PRELIMINARIES

In this section, we introduce some fundamental preliminaries that will be used throughout the paper. We first offer the definition of the notations. Then, we review the calculation process of AP and the form of AP loss. Finally, we present the cautions of AP optimisation.

3.1 | Notations

Let $\Omega = \{(\mathbf{x}_j, y_j)\}_{j=1}^N$ be the retrieval set in the retrieval systems, where (\mathbf{x}_j, y_j) indicates *j*th image and its corresponding relevant label $y_j \in \{1, 0\}$ related to the query. Note that there are a total of M queries contained in Ω , that is, $\mathcal{Q} = \{\mathbf{q}_i\}_{i=1}^M \subseteq \Omega$. For each query \mathbf{q}_i , the label y_j is assigned a relevant value, that is, '1' if \mathbf{x}_j is relevant to \mathbf{q}_i , and otherwise '0'. Thus, the retrieval set Ω can be split into the relevant \mathcal{P}_i and irrelevant \mathcal{N}_i sets, that is, $\Omega = \mathcal{P}_i \cup \mathcal{N}_i$, which are comprised by all the same class and different classes respectively. In detail, with respect to $\mathbf{q}_i, \mathcal{P}_i = \{\mathbf{x}_j \in \Omega | y_{j|\mathbf{q}_i} = 1\}$ and $\mathcal{N}_i = \{\mathbf{x}_j \in \Omega | y_{j|\mathbf{q}_i} = 0\}$. Note that each element \mathbf{x}_j in Ω is mapped to a vectorial embedding $\mathbf{v}_i \in \mathbb{R}^d$, where d is the

embedding size. In order to map all elements into the embedding space, we employ a DNN with weights θ , that is, $f_{\theta}(\cdot)$. Hence, we use $\mathbf{V} = \left\{ \mathbf{v}_{j} \in \mathbb{R}^{d} | \mathbf{v}_{j} = f_{\theta}(\mathbf{x}_{j}) \right\}_{j=1}^{N}$ to denote the vectorial embedding set.

3.2 | Average precision

AP is one of the most commonly used evaluation metrics for information retrieval tasks. It is a value of the area under the precision-recall curve. For a query \mathbf{q}_i , the set S_{Ω} of relevance scores of all elements in Ω are calculated by a selected similarity measure. In this paper, we adopt the cosine similarity, so S_{Ω} can be defined as follows:

$$S_{\Omega} = \left\{ s(\mathbf{q}_i, \mathbf{x}_j) | s(\mathbf{q}_i, \mathbf{x}_j) = \left\langle \frac{\mathbf{v}_{q_i}}{\|\mathbf{v}_{q_i}\|} \cdot \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|} \right\rangle \right\}_{j=1}^N$$
(1)

where \mathbf{v}_{q_i} and \mathbf{v}_j are the vectorical embeddings of the query \mathbf{q}_i and the element \mathbf{x}_j in Ω , $S_{\Omega} = S_{\mathcal{P}_i} \cup S_{\mathcal{N}_i}$, and $S_{\mathcal{P}_i} = \{s_k, \forall k \in \mathcal{P}_i\}$, where \mathbf{x}_j can be either an element of the relevant set \mathcal{P}_i or of the irrelevant set \mathcal{N}_i . Note that \mathbf{x}_k is always an element of the relevant set (e.g. $\mathbf{x}_k \in \mathcal{P}_i$), that is, \mathbf{q}_i and \mathbf{x}_k are of the same class. The relevant label transformation transfers the relevant labels y_j to the corresponding pairwise relevant form:

$$\forall j,k, \quad \gamma_{j,k} = 1 \tag{4}$$

where 1 is the indicator function, $\gamma_{j,k} = 1$ if $y_j = 1$, $y_k = 1$ and $\gamma_{j,k} = 0$ otherwise. Then, the ranking position of an element \mathbf{x}_j before \mathbf{x}_k in Ω can be defined as follows:

$$Rank(k, \mathcal{P}_i) = 1 + \sum_{j \in \mathcal{P}_i, j \neq k} \mathbb{1}\left\{\delta_{j,k} > 0\right\}$$
$$Rank(k, \mathcal{N}_i) = \sum_{j \in \mathcal{N}_i, j \neq k} \mathbb{1}\left\{\delta_{j,k} > 0\right\}$$
(5)

where $Rank(k, \mathcal{P}_i)$ and $Rank(k, \mathcal{N}_i)$ represent the ranking positing of the element \mathbf{x}_j in \mathcal{P}_i and Ω , and $Rank(k, \Omega) = Rank(k, \mathcal{P}_i) + Rank(k, \mathcal{N}_i)$. The AP of query \mathbf{x}_i can be calculated as follows:

$$AP_{i} = \frac{1}{|\mathcal{P}_{i}|} \sum_{k \in \mathcal{P}_{i}} precision(k) = \frac{1}{|\mathcal{P}_{i}|} \sum_{j \in \mathcal{P}_{i}} \frac{Rank(k, \mathcal{P}_{i})}{Rank(k, \Omega)}$$
$$= \frac{1}{|\mathcal{P}_{i}|} \sum_{k \in \mathcal{P}_{i}} \frac{Rank(k, \mathcal{P}_{i})}{Rank(k, \mathcal{P}_{i}) + Rank(k, \mathcal{N}_{i})}$$
$$= \frac{1}{|\mathcal{P}_{i}|} \sum_{k \in \mathcal{P}_{i}} \frac{1 + \sum_{j \in \mathcal{P}_{i}, j \neq k} \mathbb{I}\{\delta_{j,k} > 0\}}{1 + \sum_{j \in \mathcal{P}_{i}, j \neq k} \mathbb{I}\{\delta_{j,k} > 0\}}$$
(6)

 $S_{N_i} = \{s_k, \forall k \in N_i\}$ are the relevant and irrelevant score sets respectively. Each query \mathbf{q}_i thus corresponds to the relevant score set $S_{\mathcal{P}_i}$ and irrelevant score set $S_{\mathcal{N}_i}$. Thus, we can get a ranking list of images

$$\operatorname{RList}^{N}(\mathbf{q}_{i}, f_{\theta}) = \left\{ \mathbf{x}_{1}, \mathbf{x}_{2}, \cdots \mathbf{x}_{j}, \cdots \mathbf{x}_{N} | \mathbf{x}_{j} \in \Omega \right\}$$
(2)

ordered by their similarities S_{Ω} to query \mathbf{q}_i , where N is the number of returned images and f_{θ} projects \mathbf{x}_j to the feature space as vectorial embedding \mathbf{v}_j .

Next, some transformations are required to be formalised in ranking-based loss. Firstly, the difference transformation transfers the relevant score $s(\mathbf{q}_i, \mathbf{x}_j)$ to the difference form with respect to the query \mathbf{q}_i :

$$\forall j, k, \quad \delta_{j,k} = s(\mathbf{q}_i, \mathbf{x}_j) - s(\mathbf{q}_i, \mathbf{x}_k) \tag{3}$$

3.3 | AP loss

Finally, we can formulate the AP loss \mathcal{L}_{AP} by averaging all over queries, that is, the optimisation of the ranking problem can be formed as follows:

$$\min_{\theta} \mathcal{L}_{AP}(\theta) = 1 - AP(\mathbb{1}; \theta) = 1 - \frac{1}{M} \sum_{i=1}^{M} AP_i(\mathbb{1}; \theta) \quad (7)$$

where θ is the wights of the DNN and 1 corresponds to the indicator function. Note that $AP(1;\theta)$ is non-differentiable with respect to θ due to the indicator function 1, which often uses the Heaviside step function:

$$\mathcal{H}(z) = \begin{cases} 1 & z \ge 0\\ 0 & z < 0 \end{cases} \tag{8}$$

where its curve is almost all horizontal with zero or undefined gradient. Therefore, the specific smooth differentiable rank approximations need to be designed such that the standard gradient descent method can be used for optimisation.

3.4 | Cautions for AP optimisation

In order to achieve effective and better AP optimisation, some cautions for the AP itself are noteworthy. Below, we describe and analyse the key cautions for AP optimisation.

Non-differentiability. It is extremely difficult for AP optimisation due to the presence of the indicator function 1 in Equations (6) and (7). Specifically, It is the indicator function 1 that leads to the nature of non-convexity and non-differentiability in AP loss. Therefore, It is a critical challenge to design a solution that is differentiable, effective, and alternative to the indicator function 1 in AP optimisation. Meanwhile, it is currently still an open problem.

Non-decomposability. Using AP to evaluate the model performance must be in the whole test set. Therefore, AP essentially cannot linearly decompose into several batches. Limited by GPU memory resources, the current direct AP optimisation methods are based on mini-batch approximations. Consequently, this small min-batch approximation will result in a large deviation between the estimated mAP and the true mAP, that is, the above referred to as the *decomposability gap* (detailed presentation shown in Figure 2).

Limited optimisation by AP definition. AP is a metric that belongs to the ranking-based metrics, which is brittle due to the high probability of ties happening (shown in Figure 3). The main reason for this limitation lies in that the size of existing datasets is so large that there is a high degree of interclass similarity. In addition, once the positive sample score is higher than all the negative sample scores, the AP score does not change regardless of the difference between the positive and negative sample scores. In that case, if the difference between positive and negative sample scores is very slight, this case is not further optimised, which would make it highly sensitive on the test set due to the existence of potential shifts.



FIGURE 2 Illustration of estimated mAP versus batch size, including the corresponding mean and standard deviation. Intuitively, it can be noticed that the smaller the batch size is, the larger the estimated mAP (i.e. corresponding AP optimisation) bias is. mAP, mean average precision.

Therefore, the avoidance of ties has attracted attention in some works [17, 18, 23]. Finally, AP only considers both positive and negative samples before each positive sample without involving the samples after it, which may lead to local optimisation due to insufficient utilisation of sample information. Therefore, it is necessary to make full use of all samples' information for global optimisation.

4 | METHODOLOGY

In this section, we introduce the optimisation solution for the previously mentioned problems of AP-based methods. Additionally, we provide an analysis of the loss functions associated with our proposed OSAP-Loss, such as SmoothAP [25] and ROADMAP [85]. The framework of our method is illustrated in Figure 4.

4.1 | Superior ranking function

In Equation (7), it contains a discontinuous ranking function (i.e. the indicator function 1) that is non-differentiable. To address this issue, we design a novel SRF, which provides a differentiable approximation and guarantees a tighter upper bound of the AP loss. The SRF ensures robust training directly optimised using standard gradient descent methods. Specifically, we use different operations for $Rank(k, \mathcal{P}_i)$ and $Rank(k, \mathcal{N}_i)$ in Equation (6), by defining two functions $\mathcal{H}_+(\cdot)$ and $\mathcal{H}_-(\cdot)$.

For $Rank(k, \mathcal{P}_i)$, similar to ref. [85], $\mathcal{H}_+(\cdot)$ remains to be the Heaviside step function, that is, $\mathcal{H}_+(z) = \mathcal{H}(z)$ (see Figure 5a). The purpose of this operation is to avoid contradictory gradient flow for some positive samples.

For $Rank(k, N_i)$, we define a simple yet effective surrogate $\mathcal{H}_{-}(z)$ (see Figure 5b) as follows:





Ranking sample scores

FIGURE 3 Ranking sample scores may collapse when many ties happen due to the positive and negative samples with close scores. Green and red circles indicate positive and negative samples respectively.

(1): Extract features of all images 2: Forward & backward **Optimising** \mathcal{L}_{OSP} in Eq. 14 Query q_i $v_a \in \mathbb{R}^d$ ResNet-50 Pooling N×N scores N features N training images $\begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdots \end{bmatrix} N$ 1 Shared $c_1 c_2 c_3 \cdots$ c_N Gallery O Ω V ResNet-50 Pooling $\min_{\theta} \mathcal{L}_{AP}(\theta) = 1 - AP(\mathbb{1}; \theta) = 1 - \frac{1}{M} \sum_{i=1}^{M} AP_i(\mathbb{1}; \theta)$ 3: Forward & backward + gradient accumulation

FIGURE 4 The framework of our proposed OSAP-Loss employs the commonly used ResNet-50 [92] as the backbone. During the training process, N training images include an image as the query q_i and the rest images as the gallery set Ω . The optimisation procedure consists of three steps. For the first step, the network extracts the features of all images and then discards the intermediate tensors in the memory. For the second step, we first calculate the similarity matrix S_{Ω} (Equation 1), calculate the OSAP-Loss \mathcal{L}_{OSAP} containing two parts: $\min_{\theta} \mathcal{L}_{AP}(\theta)$ and \mathcal{L}_{OSP} , finally calculate the gradient of the loss regarding the features. For the last step, in order to continue backpropagation through the network, the network extracts the features from an image again, saving the intermediate tensors this time. Before the network weights are eventually updated, the gradients are accumulated, one image at a time. Best viewed in colour.



FIGURE 5 Illustration of proposed superior ranking function for substituting the Heaviside step function: (a) $\mathcal{H}_+(z)$ and (b) $\mathcal{H}_-(z)$. In this case, \mathcal{L}_{AP}^{SRF} becomes an upper bound of \mathcal{L}_{AP} , and $\mathcal{H}_-(z)$ ensures the correct gradient flow as well as avoiding gradient vanishing.

where γ and δ are hyperparameters, and γ is a scale factor and δ is an offset to make \mathcal{L}_{AP} have an upper bound. It is worth noting that $\mathcal{H}_{-}(z)$ also avoids the vanishing gradient problem for negative samples after positive samples (see the left graph in Figure 6) According to Equation (9), can achieve a superior

approximation $Rank(k, \mathcal{N}_i)_{SRF} = \sum_{j \in \mathcal{N}_i, j \neq k} \mathcal{H}_{-}(\delta_{j,k})$, obtaining the superior AP loss approximation as follows:

$$\mathcal{L}_{AP}^{SRF} = 1 - \frac{1}{M} \sum_{i=1}^{M} \frac{1}{|\mathcal{P}_i|} \sum_{k \in \mathcal{P}_i} \frac{Rank(k, \mathcal{P}_i)}{Rank(k, \mathcal{P}_i) + Rank(k, \mathcal{N}_i)_{SRF}}$$
(10)

4.2 | Memory-free mechanism

In Equation (7), AP is linearly decomposable with respect to queries \mathbf{q}_i , while AP_i is non-decomposable with respect to samples. For a query \mathbf{q}_i , the *decomposability gap* (i.e. DG_{AP}) is defined in ref. [85] as follows:

$$DG_{AP}(\theta) = \frac{1}{K} \sum_{b=1}^{K} AP_i^b(1;\theta) - AP_i(\theta)$$
(11)

where K = N/M, N is the size of the retrieval set Ω and M is the number of the queries Q. To thoroughly address the



FIGURE 6 Left: Comparison of gradient flow in AP loss for different samples in the left toy example. In contrast to SmoothAP [25] and SupAP [85], $\mathcal{L}_{OSAP} > \mathcal{L}_{AP}$ and correct gradient flow for all samples. Right: Comparison of approaches used to alleviate the decomposability gap in AP optimisation in the right toy example. Previous solutions are conducted by increasing the size of a batch with brute force, which has achieved some success. However, this method often requires a lot of computing costs. In each min-batch, \mathcal{L}_{calibr} [85] solves this problem by dictating that the positive samples' scores be greater than α and the negative samples' scores are lower than β . On the contrary, Our \mathcal{L}_{OSAP}^{MFM} employs a margin of *m* to make it more closely resemble its positive set than its negative set. *Best viewed in colour.* AP, average precision.

decomposability gap, we introduce a GPU MFM by sidestepping the memory constraints of GPU following ref. [20]. The MFM allows the model to train with arbitrary resolution images and arbitrary batch size, which makes it possible to achieve optimal AP loss, and also means that DG_{AP} can be thoroughly solved. The MFM is illustrated in Figure 4, and involves three main steps.

For the first step, the network extracts the features of all images and then discards the intermediate tensors in the memory. In the second step, we first calculate the similarity matrix S_{Ω} (Equation 1), calculate the OSAP-Loss \mathcal{L}_{OSAP} containing two parts: $\min_{\theta} \mathcal{L}_{AP}(\theta)$ and \mathcal{L}_{OSP} , finally calculate the gradient of the loss regarding the features, that is, $\frac{L_{AP}}{v_i}$. That is, in this step, we stop the network parameters updating by disrupting the back-propagation. Since a certain amount of feature vectors and similarity score matrices have to be saved, this step consumes a considerable amount of memory. During the last step, in order to continue back-propagation through the network, the network extracts the features from an image again, saving the intermediate tensors this time. Since this step consumes a lot of memory, we perform this operation on an image by image. Before the network weights are eventually updated, the gradients are accumulated in multiple steps, one image at a time.

Although the MFM sidesteps the constraint of the GPU memory resources to enable the training with a large batch size, it will contain an important overhead in memory. Especially, the last step of MFM is memory consuming, since the intermediate tensors after extracting the image features need to be stored. Therefore, this situation forces us to use the multistage backpropagation in the last step of MFM, that is, to forward one image at a time and then to back-propagate the corresponding gradient. This multistage back-propagation reduces the training speed of the network, but the utilisation of MFM decreases the number of the network iterations before the performance converge [20], and more importantly, brings a significant performance improvement (see Section 5.4). This paper focuses on performance optimisation. Since the total network training time does not increase much, the training time has not been explored accordingly in the experiments.

4.3 | Optimising samples after positive ones

Limited by AP definition, existing AP optimisation methods suffer from several problems leading to sub-optimisation (detailed discussion and analysis in Section 3.4). To address the brittleness around ties, we a score shift in order to keep a margin between the scores of different samples, which is similar to the triplet loss [68] and the circle loss [55]. Specifically, we apply a negative shift to the positive scores and a positive shift to the negative scores as follows:

$$\overrightarrow{s} \left(\mathbf{q}_i, \mathbf{x}_j \right) = \begin{cases} s \left(\mathbf{q}_i, \mathbf{x}_j \right) + \frac{\epsilon}{2} & \text{if } y_j = 0 \\ s \left(\mathbf{q}_i, \mathbf{x}_j \right) - \frac{\epsilon}{2} & \text{if } y_j = 1 \end{cases}$$
 (12)

where ϵ is a tiny margin. Then, for the query \mathbf{q}_i , we use $\delta_{j,k}$ instead of $\delta_{j,k}$ in Equation (3) to perform the difference transformation.

$$\forall j, k, \quad \stackrel{\leftrightarrow}{\delta}_{j,k} = \stackrel{\leftrightarrow}{s} \left(\mathbf{q}_i, \mathbf{x}_j \right) - \stackrel{\leftrightarrow}{s} \left(\mathbf{q}_i, \mathbf{x}_k \right)$$
(13)

Then, we consider all samples after each positive sample to fully exploit sample information for global optimisation. In addition, we also take into account that different samples should be assigned different gradients, thus allowing for high flexibility of optimisation. Therefore, we design a novel loss function to improve the optimisation of AP loss, that is, our proposed OSP loss.

4.4 | OSAP-loss

Finally, our proposed OSAP-Loss (i.e. \mathcal{L}_{OSAP}) is linearly combined with \mathcal{L}_{AP}^{SRF} and \mathcal{L}_{OSP} via a relative weight ω .

$$\mathcal{L}_{\text{OSAP}} = (1 - \omega)\mathcal{L}_{\text{AP}}^{\text{SRF}} + \omega\mathcal{L}_{\text{OSP}}$$
(16)

$$\mathcal{L}_{\text{OSP}} = \log\left(1 + \exp\frac{1}{|P|} \sum_{i}^{|P|} \sum_{j,j\neq i}^{|\Omega|} T\left\{y_{ij} \cdot w_{ij}^{p} \left[s_{i} - s_{j}\right]_{+} + \left(1 - y_{ij}\right) \cdot w_{ij}^{n} \left[\left(s_{j} + m\right) - s_{i}\right]_{+}\right\}\right)$$
(14)

$$w_{ij}^{p} = \frac{\exp([s_{i} - s_{j}]_{+})}{\sum_{y_{i}=y_{j}=1}\exp([s_{i} - s_{j}]_{+})}, w_{ij}^{n}$$
$$= \frac{\exp([(s_{j} + m) - s_{i}]_{+})}{\sum_{y_{i}=1, y_{j}=0}\exp([(s_{j} + m) - s_{i}]_{+})}$$
(15)

where $y_{ij} = 1$ if $y_i = y_j = 1$ and $y_{ij} = 0$ otherwise, *m* is a score margin. *T* is the temperature parameter that controls the scale of the sum of the score differences of all samples with respect to the query. $[\cdot]_+$ represents the non-negative operation. w_{ij}^p and w_{ij}^n inherit the advantage of relative score optimisation between positive and negative samples, which does not introduce any hyperparameters. And \mathcal{L}_{OSP} has three hyperparameters, that is, the tiny margin ϵ , the score margin *m* between positive and negative samples, and the temperature parameter *T*.

Algorithm 1 Training with OSAP-Loss

```
1: procedure TRAIN-OSAP(\Omega, C, M, n)
 2:
          \Omega: training images
 3:
          C : class labels
 4:
          M: mini-batch size
          n : number of images per class in mini-batch
 5:
 6:
 7:
          \theta \leftarrow \text{initialized}
                                     ▷ use ImageNet pre-trained parameters
 8:
          for iteration \in [1, \cdots, \text{number-of-iterations}] do
 ٩.
              loss \leftarrow 0
                                                          \triangleright set batch loss to zero
               B \leftarrow \text{Bath-sampler}(\Omega, C, M, n)
10:
11:
               for (\mathbf{q}_i, \mathbf{x}_i) \in B \times B do
12:
                   compute s(\mathbf{q}_i, \mathbf{x}_j)
13:
               end for
               for \mathbf{q}_i \in B do \triangleright use each image in the batch as query
14:
15:
                    loss \leftarrow loss + \mathcal{L}_{OSAP}
                                                                      ▷ OSAP-Loss
16:
               end for
              \theta \leftarrow \text{MINIMIZE}(\frac{loss}{|B|})
17:
                                                                    ▷ Adam update
18:
          end for
19: end procedure
```

To learn a deep model for image retrieval, we implement the OSAP-Loss based on SGD and min-batch. The training procedure using OSAP-Loss is given in Algorithm 1. Furthermore, we have added the MFM to enable further improvement of \mathcal{L}_{OSAP} , that is, \mathcal{L}_{OSAP}^{MFM} .

4.5 | Further discussion

Finally, we here compare with SmoothAP [25] and ROAD-MAP [85] (\mathcal{L}_{SupAP} and \mathcal{L}_{calibr} are two components of ROADMAP), which were recently proposed and are also most relevant to our proposed \mathcal{L}_{OSAP}^{MFM} . As shown in Figure 6, we have compared and analysed them accordingly, and there are three major differences: (1) it can well solve the issue of gradient vanishing; (2) it thoroughly solves the issue of non-decomposability in AP optimisation; (3) it solves the brittleness around ties in AP optimisation procedure, as well as enhances the separability of positive and negative samples.

Firstly, due to the sigmoid function, $\mathcal{L}_{\text{SmoothAP}}$ will suffer from gradient vanishing and contradictory gradient flow for positive samples. Although the recently proposed \mathcal{L}_{SupAP} in ROADMAP addresses the contradictory gradient flow, it does not sufficiently tackle the gradient vanishing due to the use of the sigmoid function (see the left graph of Figure 6). Instead of continuing to optimise the sigmoid function, our designed SRF $\mathcal{H}_{-}(z)$ in \mathcal{L}_{AP}^{SRF} is optimised by using improved softsign function, which can eliminate the gradient vanishing more effectively while solving the contradictory gradient flow. Secondly, it has been demonstrated that large batch sizes will obtain better performance in refs. [17, 20, 23, 25, 85]. The reason lies in the fact that the AP is not linearly decomposed into multiple batches for optimisation. Therefore, the larger the batch size, the closer the estimated AP value is to the true AP, and the smaller the DG_{AP} is. To overcome the memory constraints of GPU, our proposed MFM is to train the deep model with a large batch size, which can thoroughly solve the nondecomposability of AP (see the right graph of Figure 6).

Thirdly, for ranking-based metrics, one of the major drawbacks is the easy occurrence of ties, especially in large datasets. Our \mathcal{L}_{OSP} solves the brittleness around ties with score shift operation as well as improves the robustness of the learned model. These differences provide a significant performance improvement over several image retrieval datasets (Section 5).

Furthermore, we discuss the comparison of our method with other methods in terms of computational complexity. In total, we compare three types of losses, that is, proxy-based, pairbased and AP-based losses in Table 1. For proxy-based and pairbased losses, we can find that their computational complexity is much larger compared to AP-based losses. Meanwhile, in contrast to the general ranking function, they are all optimised based on local samples or sample pairs, which makes them easy to fall into local optimisation. For QS-Suitable [83] and SoftBin* [20], although they use AP as the optimisation objective and have a relatively small computational complexity, they do not employ a general ranking function, which leads to a very coarse AP approximation. Therefore, it is necessary to use an appropriate general ranking function for AP optimisation. With the help of the general ranking function, the AP-based losses can be continuously optimised towards the global optimum. However, these methods (BlackBox [24], PNP-D_q [84], FastAP [17], and SoDeep [22]) perform a coarse AP approximation even by operating a general ranking function. Aiming to design tighter

smooth AP approximation (SmoothAP [25] and SupAP [25]), our \mathcal{L}_{AP}^{SRF} proposed to design smooth differentiable approximation of the ranking function, where the core idea is similar to SupAP [85] with the same computational complexity $(O((\mathcal{N} + \mathcal{P})\mathcal{P}))$. We can see that the computational complexity of SmoothAP [25] is $O(\mathcal{M}^2)$, which is larger than that of SupAP [25] and our \mathcal{L}_{AP}^{SRF} . Compared to these methods, it is noted that our method is superior to both of them, as already discussed and analysed above.

5 | EXPERIMENTS

In this section, we first introduce the implementation details of our proposed method, the test protocol, and the benchmarks used for experiments, including three remote sensing image datasets and three CIR datasets. Then, we provide the relevant ablation study results and give a detailed analysis. Finally, we compare OSAP-Loss with several state-of-the-art methods.

5.1 | Datasets

We evaluate OSAP-Loss on the following three RSIR datasets (UCMD [40], NWPUD [41], and PatternNet [1]) and three

Loss	Туре	Computational complexity	General ranking (AP approximation)
Proxy-Anchor [80]	Proxy	$O(\mathcal{MC})$	X (-)
Proxy-NCA [70]		$O(\mathcal{MC})$	X (-)
SoftTriple [76]		$Oig(\mathcal{MCU}^2ig)$	X (-)
Contrastive [93, 94]	Pair	$O(\mathcal{M}^2)$	X (-)
Triplet (smart) [74]		$Oig(\mathcal{M}^2ig)$	X (-)
Triplet (semi-hard) [68]		$Oig(\mathcal{M}^3/\mathcal{B}^2)$	X (-)
<i>N</i> -pair [69]		$O(\mathcal{M}^3)$	X (-)
Lifted structure [43]		$O(\mathcal{M}^3)$	X (-)
QS-suitable [83]	AP	$O(\mathcal{N} \log \mathcal{P})$	X (Coarse)
SoftBin* [20]		$O(\mathcal{NP})$	X (Coarse)
BlackBox [24]		$O(\mathcal{N} \mathrm{log}\mathcal{N})$	✔(Coarse)
PNP-D _q [84]		$O(\mathcal{NP})$	✔(Coarse)
FastAP [17]		$O((\mathcal{N}+\mathcal{P})\mathcal{L})$	✔(Coarse)
SoDeep [22]		$Oig((\mathcal{N}+\mathcal{P})\mathcal{H}^2ig)$	✔(Coarse)
SmoothAP [25]		$Oig(\mathcal{M}^2ig)$	✔(Smooth)
SupAP [85]		$O((\mathcal{N}+\mathcal{P})\mathcal{P})$	✔(Smooth)
OSAP $\left(\mathcal{L}_{AP}^{SRF}\right)$		$O((\mathcal{N}+\mathcal{P})\mathcal{P})$	✔(Smooth)

Note: For proxy based and pair based losses, the numbers of training samples, classes, batch size in each epoch, and proxies of each class are represented by $\mathcal{M}, \mathcal{C}, \mathcal{B}$, and \mathcal{U} respectively. For AP based losses, \mathcal{P} and \mathcal{N} ($\mathcal{P} + \mathcal{N} = \mathcal{M}$) denote the number of relevant (positive) and irrelevant (negative) samples respectively. The number of bins is denoted by \mathcal{L} for FastAP. For SoDeep, \mathcal{H} denotes the hidden state size ($\mathcal{H} \approx \mathcal{N}$) in LSTM. SmoothAP, SupAP, and our \mathcal{L}_{AP}^{SRF} are smooth AP approximations, which makes the AP calculation more accurate in practice.

Abbreviation: AP, average precision.

TABLE 1 Comparison of the computational complexity for three types of losses, including proxy-based, pair-based and AP-based loss functions

CIR datasets (CUB-200-2011 [42], SOP [43], and INaturalist-2018 [44]). The number of images, classes, and the average number of images per class are statistics in Table 2. We have counted the number of images (# Imges), classes (# Classes), and images/class (# Avg) in the training set and testing set for each dataset. Next, we describe the datasets in more detail.

- 1) RSIR Datasets
 - **UCMD** [40] is collected by the University of California Merced from the United States Geological Survey (USGS) for land use and cover. This dataset has 2100 high-resolution images with 21 different scene classes, which are 256 × 256 pixels and are 0.3 m in spatial resolution. We follow the standard protocol [95], which uses 50% images of each class as the training set and the

rest as the test set. The sample images of each class are shown in Figure 7.

- NWPU-RESISC45 (NWPUD) [41] is collected by the North-western Polytechnical University from Google Earth. This dataset contains 31,500 images of 45 classes. Each image is also 256 × 256 pixels and varies from 0.2 to 30 m in spatial resolution. All images cover more than 100 countries. Following the standard data splitting protocol [41], 80% images for training and the rest of 20% images for evaluation. Figure 8 illustrates samples of each class.
- PatternNet [1] is a large-scale high-resolution dataset collected from Google Earth imagery or via Google Map API for RSIR. This dataset includes 30,400 images with 38 classes, and there are 800 images of 256 × 256 pixels

Dataset		# Images	# Classes	# Avg
Remote sensing image retrieval datasets	UCMD train [40]	1100	21	52.4
	UCMD test [40]	1100	21	52.4
	NWPUD train [41]	25,200	45	560.0
	NWPUD test [41]	6300	45	140.0
	PatternNet train [1]	24,320	38	640.0
	PatternNet test [1]	6080	38	160.0
Conventional image retrieval datasets	CUB train [42]	5864	100	58.6
	CUB test [42]	5924	100	59.2
	SOP train [43]	59,551	11,318	5.3
	SOP test [43]	60,502	11,316	5.3
	INaturalist train [44]	325,846	5690	57.3
	INaturalist test [44]	136,093	2452	55.5

TABLE 2 Dataset composition for training and evaluation



FIGURE 7 Illustration of sample images of different scenes from UCMD dataset.



FIGURE 8 Illustration of sample images of different scenes from NWPUD dataset.

per class. The spatial resolution ranges from 0.062 to 4.693 m. We follow the standard splitting protocol in ref. [1], that is, 80% training set and 20% testing set. Figure 9 shows corresponding sample images in the dataset.

- 2) CIR Datasets
 - CUB-200-2011 [42] is a challenging dataset of 200 bird species. This dataset has 11,788 images. Following the standard protocol, we use the first 100 classes with 5864 images as the training set and the rest of the 100 classes with 5924 images as the testing set.
 - Stanford Online Product (SOP) is an online products dataset collected by Oh Song et al. [43] from eBay.com. This dataset contains 120,053 images with 22,634 classes. Referring to ref. [43], we use 11,318 classes with 59,551 images as the training set and 11,316 classes with 60,502 images as the testing set.
 - **INaturalist-2018** [44] is a real-world dataset with 461,939 images. This dataset includes 8142 iNaturalist species. Similar to ref. [25], we use 325,846 images of 5690 classes for training and 136,093 images of 2452 classes for testing.

5.2 | Implementation details

Our method includes training and evaluation, both implemented in the PyTorch framework. In our experiments, we use ResNet-50 [92] as the backbone. The model is initialized by the pre-trained parameters on ImageNet [96]. We use Adam [97] to optimize model training, and use standard data pre-processing procedure and augmentation strategy during training, that is, normalizing all images to 256×256 and randomly cropping them to 224×224 as input images. The crops are flipped horizontally with a 50% probability. During the evaluation, all images are normalized to 256×256 and then center cropped to 224 \times 224. We use 333 as the fixed random seed in our all experiments to avoid seed performance fluctuations. We set the buffer size to 512 in MFM. For RSIR datasets, we train the models with a learning rate of 10^{-6} for 400 epochs. For CIR datasets, we also set the learning rate to 10^{-6} on CUB for 200 epochs. For SOP and INaturalist, we set the initial learning rate to 10^{-5} and drop the learning rate by 70% on 30 and 70 epochs. And the models are trained on SOP for 100 epochs and on INaturalist for 90 epochs. For all experiments in Section 5.4 and the models in Section 5.5, we use $\omega = 0.2$ for

13



FIGURE 9 Illustration of sample images of different scenes from PatternNet dataset.

 $\mathcal{L}_{\text{OSAP}}$ in Equation (16), $\gamma = 0.001$ and $\delta = 0.65$ for $\mathcal{L}_{\text{AP}}^{\text{SRF}}$ in Equation (10), $\epsilon = 0.20$ on RSIR datasets, $\epsilon = 0.02$ on CIR datasets, and m = 0.25 for \mathcal{L}_{OSP} in Equation (14).

5.3 | Test protocol

To evaluate the model performance, we describe here the used protocols in our experiments. For all image retrieval datasets, each sample from each class is utilised as the query \mathfrak{q}_i in turn, and the retrieval set Ω is made up of all the remaining samples. For RSIR, we use standard precision at k (P@k, precision of the top-k retrieval results) and mean average precision (mAP, detailed calculation given in Section 3). For CIR, we use standard recall at k (R@k, hit rate of the top-k retrieval results) and mAP at R (mAP@R) [79]. Below, we introduce the computation of the mentioned evaluation metrics. **Precision@K** The Precision@K in Equation (17) is often used in the retrieval task. For each query, the Precision@K is precision at k, that is, precision of the top-k retrieval results. The Precision@K is averaged over all the queries.

$$P@K = \frac{\sum_{i=1}^{M} P_k(i)}{M}, \text{ where } P_k(i) \text{ is the precision at } k. \quad (17)$$

Recall@K The Recall@K in Equation (18) is another evaluation metric for the retrieval task. For each query, the Recall@K is Recall at k, Recall@K = 1 if a positive sample appears in the top-k retrieval results and Recall@K = 0 otherwise. The Recall@K is averaged over all the queries in Equation (18), where R(i) = 1 if a positive sample has a higher ranking than K, and R(i) = 0 otherwise.

$$R@K = \frac{\sum_{i=1}^{M} R(i)}{M}$$
(18)

mAP@R The mAP@R in Equation (19) is proposed in ref. [79]. It is less noisy and easily obtains the better performance of the model. The mAP@R is essentially a partial AP, which is an AP of R retrieved positive samples with respect to a query.

$$mAP@R = \frac{\sum_{j=1}^{R} P(j)}{R}, \text{ where } P(j) \text{ is the precision at } j.$$
(19)

5.4 | Ablation study

For a fair comparison, we conduct all our experiments here with the same settings (Backbone: ResNet-50; Batch size: 64; Embedding size: 512). Such comparisons allow us to directly observe the impact of different parts or parameters on the final performance.

2682322, 0, Downloaded from https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cit2.12151, Wiley Online Library on [2505/2023]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

5.4.1 | Impact of our designed components

To investigate more deeply the impact of our designed components, the relevant results of the ablation experiments are shown in Tables 3 and 4. Here, Table 3 shows the results on RSIR datasets and Table 4 shows the results on CIR datasets. In Table 3, it is clear to see that the performance of the SmoothAP baseline [25] is already impressive on three RSIR datasets. Meanwhile, we find that the performance of using \mathcal{L}_{AP}^{SRF} or \mathcal{L}_{OSP} alone does not perform better or slightly degrades on UCMD, NWPUD, and SOP. We presume that such performance degradation may be caused by local optimisation arising from the low diversity caused by the small number of samples within the batch. The performance improvement from using MFM alone and from using MFM on \mathcal{L}_{OSAP} shows that our presumption is plausible. In Table 4, it can be obviously observed that compared to SmoothAP [25], our method performs a great improvement by replacing the sigmoid with \mathcal{H}_+ and \mathcal{H}_{-} for \mathcal{L}_{AP}^{SRF} in Equation (10) especially on CUB. Note that the use of \mathcal{L}_{OSP} in Equation (14) on \mathcal{L}_{AP}^{SRF} further boosts the experimental performance. After using MFM, we can see that \mathcal{L}_{OSAP} receives a further performance improvement of mAP/mAP@R on RSIR datasets (~3pt on UCMD, ~7pt on NWPUD, and 0.01pt on PatternNet) and CIR datasets (~6pt

Dataset	SRF	OSP	MFM	mAP	P@5	P@10	P@50	P@100	P@1000
UCMD	x	x	X	96.49	99.98	99.98	99.91	78.79	7.90
	✓	x	X	96.33	99.88	99.85	99.77	78.84	7.90
	x	✓	X	93.20	98.26	97.99	96.88	76.98	7.90
	x	x	√	96.60	99.99	99.99	99.98	78.99	7.90
	✓	✓	x	96.29	99.80	99.76	99.68	78.85	7.90
	1	✓	√	99.54	100.00	100.00	99.92	78.98	7.92
NWPUD	x	x	x	92.53	95.35	95.16	94.37	93.00	13.89
	1	x	x	90.38	94.31	93.99	92.94	91.15	13.88
	x	✓	x	88.94	92.84	92.40	90.95	89.87	13.90
	x	x	√	97.79	98.49	98.50	98.30	97.87	13.89
	1	✓	x	92.60	94.83	94.48	93.72	92.63	13.90
	1	✓	√	99.73	100.00	100.00	99.91	99.80	13.93
PatternNet	x	x	x	99.93	99.95	99.95	99.95	99.95	55.90
	✓	x	x	99.87	99.91	99.90	99.88	99.87	55.88
	x	✓	x	99.97	99.97	99.97	99.96	99.96	55.90
	x	x	1	99.98	99.98	99.98	99.98	99.98	55.90
	1	1	x	99.98	99.98	99.98	99.98	99.98	55.90
	1	1	1	99.99	100.00	100.00	99.99	99.99	56.14
	Dataset UCMD NWPUD	Dataset SRF UCMD X X X X X X X X X X X X X X X X X X X	DatasetSRFOSPUCMDXXIXXXXXXXXIIIIXXIXXIXXIXIIXIIXIII </td <td>DatasetSRFOSPMFMUCMDXXXICMDXXXICMDIIXICMIIIICMII</td> <td>DatasetSRFOSPMFMmAPUCMDXXX96.49\checkmarkXX96.33$\checkmark$$\checkmark$X93.20$\chi$$\checkmark$$\chi93.20\chi$$\checkmark$$\chi96.69\chi$$\checkmark$$\checkmark96.29\checkmark$$\checkmark$$\checkmark96.29\checkmark$$\checkmark$$\chi$99.54NWPUD$\chi$$\chi$$\chi$NWPUD$\chi$$\chi$$\chi$$\chi$$\chi$$\chi90.38\chi$$\checkmark$$\chi90.31\chi$$\chi$$\chi92.60\checkmark$$\checkmark$$\chi99.93\checkmark$$\chi$$\chi$$\chi$PatternNet$\chi$$\chi$$\chi$$\chi$$\chi$$\chi99.93\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi99.98\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$$\chi$$\chi$$\varphi$$\chi$<</td> <td>DatasetSRFOSPMFMmAPP@5UCMD$X$$X$$X$96.4999.98$\checkmark$$X$$X$96.3399.88$X$$\checkmark$$\chi$93.2098.26$X$$\checkmark$$\checkmark$96.6099.99$\checkmark$$\checkmark$$\checkmark$96.2999.80$\checkmark$$\checkmark$$\checkmark$96.2999.80$\checkmark$$\checkmark$$\checkmark$99.54100.00NWPUD$X$$X$$\checkmark$92.5395.35$\checkmark$$\chi$$\chi$$\chi$90.3894.31$X$$\checkmark$$\chi$90.3894.31$\chi$$\checkmark$$\chi$$\varphi$99.7390.84$\checkmark$$\checkmark$$\chi$$\varphi$99.73100.00PatternNet$\chi$$\chi$$\chi$99.9399.95$\checkmark$$\chi$$\chi$$\varphi$99.9799.91$\chi$$\chi$$\chi$$\varphi$99.9899.98$\checkmark$$\chi$$\checkmark$99.9899.98$\checkmark$$\chi$$\checkmark$99.9899.98$\checkmark$$\chi$$\chi$$\varphi$99.9899.98$\checkmark$$\chi$$\chi$$\varphi$99.9899.98$\checkmark$$\chi$$\chi$$\varphi$99.9899.98$\checkmark$$\chi$$\varphi$99.9899.98$\checkmark$$\psi$$\varphi$$\varphi$99.99100.00</td> <td>DatasetSRFOSPMFM$mAP$$P@5$$P@10$UCMD$X$$X$$X$96.4999.9899.98$\checkmark$$X$$X$96.3399.8899.85$X$$\checkmark$$X$93.2098.2697.99$X$$X$$\checkmark$96.6099.9999.99$\checkmark$$\checkmark$$\checkmark$96.2999.8099.76$\checkmark$$\checkmark$$\checkmark$99.54100.00100.00NWPUD$X$$X$$\chi$92.5395.3595.16$\checkmark$$\chi$$\chi$$\chi$90.3894.3193.99$\chi$$\checkmark$$\chi$99.7699.8492.40$\chi$$\chi$$\chi$$\varphi$91.3393.99$\chi$$\chi$$\chi$$\varphi$91.3393.99$\chi$$\chi$$\chi$$\varphi$91.3393.99$\chi$$\chi$$\chi$92.6094.8394.48$\checkmark$$\checkmark$92.6094.8394.48$\checkmark$$\checkmark$99.9399.9599.95$\checkmark$$\chi$$\chi$$\varphi$99.9399.91PatternNet$X$$\chi$$\varphi$99.9399.91$\chi$$\chi$$\chi$$\varphi$99.9799.97$\chi$$\chi$$\chi$$\varphi$99.9899.98$\chi$$\chi$$\chi$$\varphi$99.9899.98$\chi$$\chi$$\chi$$\varphi$99.9899.98$\chi$$\chi$$\varphi$99.9899.9899.98$\chi$<!--</td--><td>DatasetSRFOSPMFMmAPP@5P@10P@50UCMD$X$$X$$X$96.4999.9899.9899.91$\checkmark$$X$$X$96.3399.8899.8599.77$\chi$$\checkmark$$X$93.2098.2697.9996.88$\chi$$\checkmark$$\chi$96.6099.9999.9999.98$\checkmark$$\checkmark$$\checkmark$96.2999.8099.7699.68$\checkmark$$\checkmark$$\checkmark$99.54100.00100.0099.92NWPUD$\chi$$\chi$$\chi$92.5395.3595.1694.37$\checkmark$$\chi$$\chi$$\chi$90.3894.3193.9992.94$\chi$$\chi$$\chi$$\chi$92.6094.8394.4893.72$\chi$$\chi$$\chi$$\chi$99.9799.9599.9599.95$\chi$$\chi$$\chi$99.37100.00100.0099.91PatternNet$\chi$$\chi$$\chi$99.8799.9199.9099.98$\chi$$\chi$$\chi$$\varphi$99.9799.9799.9799.97$\chi$$\chi$$\chi$$\varphi$99.9799.9799.9999.98$\chi$$\chi$$\chi$$\varphi$99.9799.9799.99$\chi$$\chi$$\chi$99.9799.9799.9799.97$\chi$$\chi$$\chi$99.9799.9799.9799.98$\chi$$\chi$$\chi$99.9899.9899.9899.98</td><td>DatasetSRFOSPMFM$mAP$$P@5$$P@10$$P@50$$P@10$UCMDXXX96.4999.9899.9899.9178.79\checkmarkXX96.3399.8899.8599.7778.84$\chi$$\checkmark$X93.2098.2697.9996.8876.98$\chi$$\chi$$\checkmark$96.6099.9999.9999.9878.99$\checkmark$$\checkmark$$\chi$96.2999.8099.7699.6878.85$\checkmark$$\checkmark$90.2999.8099.7699.6878.85$\checkmark$$\checkmark$99.54100.00100.0099.9278.98NWPUD$\chi$$\chi$$\chi$92.5395.1694.3793.00$\checkmark$$\chi$$\chi$90.3894.3193.9992.9491.15$\chi$$\checkmark$$\chi$90.3894.3193.9992.9491.15$\chi$$\checkmark$$\chi$90.3894.3193.9992.9491.15$\chi$$\checkmark$$\chi$90.3894.8493.7292.63$\checkmark$$\chi$$\checkmark$90.7798.4998.5098.3097.87$\chi$$\checkmark$$\chi$99.73100.00100.0099.9199.95$\checkmark$$\checkmark$99.9399.9599.5599.9599.95$\checkmark$$\checkmark$99.73100.00100.0099.9499.96$\chi$$\chi$$\checkmark$99.9799.9799.9599.9599.95<tr< td=""></tr<></td></td>	DatasetSRFOSPMFMUCMDXXXICMDXXXICMDIIXICMIIIICMII	DatasetSRFOSPMFM mAP UCMDXXX96.49 \checkmark XX96.33 \checkmark \checkmark X93.20 χ \checkmark χ 93.20 χ \checkmark χ 96.69 χ \checkmark \checkmark 96.29 \checkmark \checkmark \checkmark 96.29 \checkmark \checkmark χ 99.54NWPUD χ χ χ NWPUD χ χ χ χ χ χ 90.38 χ \checkmark χ 90.31 χ χ χ 92.60 \checkmark \checkmark χ 99.93 \checkmark χ χ χ PatternNet χ χ χ χ χ χ 99.93 χ χ χ φ χ χ χ 99.98 χ χ χ φ χ <	DatasetSRFOSPMFMmAPP@5UCMD X X X 96.4999.98 \checkmark X X 96.3399.88 X \checkmark χ 93.2098.26 X \checkmark \checkmark 96.6099.99 \checkmark \checkmark \checkmark 96.2999.80 \checkmark \checkmark \checkmark 96.2999.80 \checkmark \checkmark \checkmark 99.54100.00NWPUD X X \checkmark 92.5395.35 \checkmark χ χ χ 90.3894.31 X \checkmark χ 90.3894.31 χ \checkmark χ φ 99.7390.84 \checkmark \checkmark χ φ 99.73100.00PatternNet χ χ χ 99.9399.95 \checkmark χ χ φ 99.9799.91 χ χ χ φ 99.9899.98 \checkmark χ \checkmark 99.9899.98 \checkmark χ \checkmark 99.9899.98 \checkmark χ χ φ 99.9899.98 \checkmark χ χ φ 99.9899.98 \checkmark χ χ φ 99.9899.98 \checkmark χ φ 99.9899.98 \checkmark ψ φ φ 99.99100.00	DatasetSRFOSPMFM mAP $P@5$ $P@10$ UCMD X X X 96.4999.9899.98 \checkmark X X 96.3399.8899.85 X \checkmark X 93.2098.2697.99 X X \checkmark 96.6099.9999.99 \checkmark \checkmark \checkmark 96.2999.8099.76 \checkmark \checkmark \checkmark 99.54100.00100.00NWPUD X X χ 92.5395.3595.16 \checkmark χ χ χ 90.3894.3193.99 χ \checkmark χ 99.7699.8492.40 χ χ χ φ 91.3393.99 χ χ χ φ 91.3393.99 χ χ χ φ 91.3393.99 χ χ χ 92.6094.8394.48 \checkmark \checkmark 92.6094.8394.48 \checkmark \checkmark 99.9399.9599.95 \checkmark χ χ φ 99.9399.91PatternNet X χ φ 99.9399.91 χ χ χ φ 99.9799.97 χ χ χ φ 99.9899.98 χ χ χ φ 99.9899.98 χ χ χ φ 99.9899.98 χ χ φ 99.9899.9899.98 χ </td <td>DatasetSRFOSPMFMmAPP@5P@10P@50UCMD$X$$X$$X$96.4999.9899.9899.91$\checkmark$$X$$X$96.3399.8899.8599.77$\chi$$\checkmark$$X$93.2098.2697.9996.88$\chi$$\checkmark$$\chi$96.6099.9999.9999.98$\checkmark$$\checkmark$$\checkmark$96.2999.8099.7699.68$\checkmark$$\checkmark$$\checkmark$99.54100.00100.0099.92NWPUD$\chi$$\chi$$\chi$92.5395.3595.1694.37$\checkmark$$\chi$$\chi$$\chi$90.3894.3193.9992.94$\chi$$\chi$$\chi$$\chi$92.6094.8394.4893.72$\chi$$\chi$$\chi$$\chi$99.9799.9599.9599.95$\chi$$\chi$$\chi$99.37100.00100.0099.91PatternNet$\chi$$\chi$$\chi$99.8799.9199.9099.98$\chi$$\chi$$\chi$$\varphi$99.9799.9799.9799.97$\chi$$\chi$$\chi$$\varphi$99.9799.9799.9999.98$\chi$$\chi$$\chi$$\varphi$99.9799.9799.99$\chi$$\chi$$\chi$99.9799.9799.9799.97$\chi$$\chi$$\chi$99.9799.9799.9799.98$\chi$$\chi$$\chi$99.9899.9899.9899.98</td> <td>DatasetSRFOSPMFM$mAP$$P@5$$P@10$$P@50$$P@10$UCMDXXX96.4999.9899.9899.9178.79\checkmarkXX96.3399.8899.8599.7778.84$\chi$$\checkmark$X93.2098.2697.9996.8876.98$\chi$$\chi$$\checkmark$96.6099.9999.9999.9878.99$\checkmark$$\checkmark$$\chi$96.2999.8099.7699.6878.85$\checkmark$$\checkmark$90.2999.8099.7699.6878.85$\checkmark$$\checkmark$99.54100.00100.0099.9278.98NWPUD$\chi$$\chi$$\chi$92.5395.1694.3793.00$\checkmark$$\chi$$\chi$90.3894.3193.9992.9491.15$\chi$$\checkmark$$\chi$90.3894.3193.9992.9491.15$\chi$$\checkmark$$\chi$90.3894.3193.9992.9491.15$\chi$$\checkmark$$\chi$90.3894.8493.7292.63$\checkmark$$\chi$$\checkmark$90.7798.4998.5098.3097.87$\chi$$\checkmark$$\chi$99.73100.00100.0099.9199.95$\checkmark$$\checkmark$99.9399.9599.5599.9599.95$\checkmark$$\checkmark$99.73100.00100.0099.9499.96$\chi$$\chi$$\checkmark$99.9799.9799.9599.9599.95<tr< td=""></tr<></td>	DatasetSRFOSPMFMmAPP@5P@10P@50UCMD X X X 96.4999.9899.9899.91 \checkmark X X 96.3399.8899.8599.77 χ \checkmark X 93.2098.2697.9996.88 χ \checkmark χ 96.6099.9999.9999.98 \checkmark \checkmark \checkmark 96.2999.8099.7699.68 \checkmark \checkmark \checkmark 99.54100.00100.0099.92NWPUD χ χ χ 92.5395.3595.1694.37 \checkmark χ χ χ 90.3894.3193.9992.94 χ χ χ χ 92.6094.8394.4893.72 χ χ χ χ 99.9799.9599.9599.95 χ χ χ 99.37100.00100.0099.91PatternNet χ χ χ 99.8799.9199.9099.98 χ χ χ φ 99.9799.9799.9799.97 χ χ χ φ 99.9799.9799.9999.98 χ χ χ φ 99.9799.9799.99 χ χ χ 99.9799.9799.9799.97 χ χ χ 99.9799.9799.9799.98 χ χ χ 99.9899.9899.9899.98	DatasetSRFOSPMFM mAP $P@5$ $P@10$ $P@50$ $P@10$ UCMDXXX96.4999.9899.9899.9178.79 \checkmark XX96.3399.8899.8599.7778.84 χ \checkmark X93.2098.2697.9996.8876.98 χ χ \checkmark 96.6099.9999.9999.9878.99 \checkmark \checkmark χ 96.2999.8099.7699.6878.85 \checkmark \checkmark 90.2999.8099.7699.6878.85 \checkmark \checkmark 99.54100.00100.0099.9278.98NWPUD χ χ χ 92.5395.1694.3793.00 \checkmark χ χ 90.3894.3193.9992.9491.15 χ \checkmark χ 90.3894.3193.9992.9491.15 χ \checkmark χ 90.3894.3193.9992.9491.15 χ \checkmark χ 90.3894.8493.7292.63 \checkmark χ \checkmark 90.7798.4998.5098.3097.87 χ \checkmark χ 99.73100.00100.0099.9199.95 \checkmark \checkmark 99.9399.9599.5599.9599.95 \checkmark \checkmark 99.73100.00100.0099.9499.96 χ χ \checkmark 99.9799.9799.9599.9599.95 <tr< td=""></tr<>

Note: The results are shown with the different designed components, and the SmoothAP baseline (i.e. SRF (X), OSP (X), and MFM (X)) on UCMD, NWPUD, and PatternNet datasets. The results in bold indicate the best performance.

Abbreviations: mAP, mean average precision; MFM, memory-free mechanism; OSAP, Optimising samples after positive ones & average precision loss; OSP, Optimising Samples after Positive ones; SRF, superior ranking function.

				SOP		CUB		INatu	ralist
Method	SRF	OSP	MFM	R@1	mAP@R	R@1	mAP@R	R@1	mAP@R
SmoothAP [25]	x	x	x	80.1	54.6	62.1	23.9	59.7	20.7
OSAP	✓	x	x	79.5	52.7	65.0	25.5	61.6	21.7
	x	✓	x	77.3	49.6	65.0	23.7	61.2	21.1
	x	x	1	82.2	56.3	66.2	26.1	66.8	26.3
	✓	✓	x	78.7	51.4	67.1	26.5	66.3	25.4
	✓	✓	✓	84.2	57.1	69.6	26.1	70.6	26.7

Note: The results are shown with the different designed components, and the SmoothAP baseline on SOP, CUB and INaturalist datasets. The results in bold indicate the best performance.

Abbreviations: mAP, mean average precision; MFM, memory-free mechanism; OSAP, Optimising samples after positive ones & average precision loss; OSP, Optimising Samples after Positive ones; SOP, Stanford Online Product; SRF, superior ranking function.

 $T\,A\,B\,L\,E\,\,5$ $\,$ Ablation study over different hyperparameters of scale factor γ

	UCMD	CUB	
γ	mAP	R@1	mAP@R
1	93.24	40.19	11.43
0.1	96.62	52.58	18.58
0.01	96.62	61.21	23.87
0.001	96.33	64.30	24.38
0.0001	94.07	52.53	13.49

Note: The results are reported on UCMD and CUB. The results in bold indicate the best performance.

on SOP, 0.6pt on CUB, and 1.3pt on INaturalist). In Table 3, \mathcal{L}_{AP}^{SRF} has a greater contribution compared to \mathcal{L}_{OSP} , and the result can be further improved by adding \mathcal{L}_{OSP} to \mathcal{L}_{AP}^{SRF} on UCMD and NWPUD. Differently, \mathcal{L}_{AP}^{SRF} and \mathcal{L}_{OSP} contribute approximately the same to the performance improvement on CUB and INaturalist in Table 4. However, \mathcal{L}_{AP}^{SRF} and \mathcal{L}_{OSP} perform not especially well on SOP. We postulate that the reason may be # Avg is too small to optimise on SOP, as shown in Table 2. Note that the performance improvement of adding MFM is much more significant. This is explained by the fact that the MFM reduces the decomposability gap in Equation (11). We can see that MFM is very effective for improving performance and brings large gains both on RSIR datasets and CIR datasets.

5.4.2 | Impact of the hyperparameters

In this section, we will exploit the impact of the hyperparameters in $\mathcal{L}_{\text{OSAP}}$ without MFM. We analyse the impact of five hyperparameters, that is, the scale factor γ and the offset δ in Equation (9), the tiny margin ϵ in Equation (12), the score margin m and the temperature parameter T in Equation (14), and the relative weight ω in Equation (16) on RSIR and CIR tasks respectively. Note that we only change one parameter at a time.



TABLE 4 Ablation study on our

designed OSAP-Loss

FIGURE 10 Illustration of the effect of offset δ on (a) UCMD and (b) CUB datasets respectively.

- The hyperparameters in SRF. Here, we first exploit the effect of different hyperparameters in L^{SRF}_{AP}, including the scale factor γ, and the offset δ. We vary the two hyperparameters in a series of values to observe the effect on the experimental results, and the results are shown in Table 5 and Figure 10.
 - ^o Effect of scale factor γ. γ governs the smoothing of softsign that is used to approximate the indicator function in \mathcal{L}_{AP}^{SRF} . The scale factor is essential in lots of AP approximation methods [25, 84, 85]. It determines the region range size for gradient backpropagation of negative samples with high scores. We set γ to different values and the performances are shown in Table 5. The ablation results show that a value of 0.001 leads to the best performance (i.e. mAP@R and R@1), which achieves the trade-off of AP approximation and gradient optimisation. Furthermore, the setting of the suitable scale factor is supposed to ensure a smooth AP approximation while also providing a large enough optimisation region for the gradient optimisation, thus achieving considerable optimisation.
 - Effect of offset δ . δ controls the upper bound of AP in \mathcal{H}_- and determines the degree of AP optimisation. We vary δ from 0 to 1 (with 0.05 as the interval) and visualise the results on UCMD and CUB in Figure 10. For UCMD, it can be observed that the change of δ' value has very little impact on the performance. For CUB, we can see that the performance is gradually improving when the value of δ from 0 to 0.65. Note that the best

performance is achieved when the value of δ is 0.65. There is a significant degradation when $\delta > 0.65$ possibly due to AP's overestimation.

- 2) The hyperparameters in OSP. Then, we exploit the effect of different hyperparameters in \mathcal{L}_{OSP} . It includes three hyperparameters, that is, the tiny margin ϵ , the score margin m, and the temperature parameter T. ϵ is set to 0.20 on UCMD and 0.02 on CUB empirically in all experiments. Therefore, we study the effect of T and m. And the results are shown in Table 6.
 - Effect of score margin *m*. *m* determines the score margin between positive and negative samples. Here, we set *m* to 0, 0.15, 0.25, 0.35, and 0.45. And we report the results in Table 6. For UCMD, the performance is relatively stable, holding at around 95.00%. For CUB, it can be observed that mAP@R is consistently stable at 23.7% in [0.25, 0.45], but R@1 increases and then decreases in [0.25, 0.45]. And R@1 and mAP@R achieve the best results (65.0% R@1 and 23.7% mAP@R) when *m* is 0.25.
 - Effect of temperature parameter *T*. *T* determines the scale of the sum of optimisation scores of positive and negative samples. It is critical in optimisation of \mathcal{L}_{OSP} . We vary *T* from 45 to 85 (with 10 as the interval) and report the results in Table 6. For UCMD, \mathcal{L}_{OSP} achieve the best performance (95.58% mAP) when T = 45. For CUB, we can see that the best performance (65.0% R@1 and 23.7% mAP@R) when T = 65. It can be observed

TABLE 6 Ablation study over different hyperparameters: score margin m and temperature parameter T

	UCMD	CUB			UCMD	CUB	CUB			
m	mAP	R@1	mAP@R	Т	mAP	R@1	mAP@R			
0	95.45	65.1	23.6	45	95.58	65.0	23.4			
0.15	95.07	64.9	23.7	55	95.30	64.9	23.6			
0.25	95.05	65.0	23.7	65	95.05	65.0	23.7			
0.35	94.98	64.6	23.7	75	94.84	64.7	23.7			
0.45	94.89	64.6	23.7	85	94.68	64.1	23.6			

Note: The results are reported on UCMD and CUB. The results in **bold** indicate the best performance.



FIGURE 11 Illustration of the effect of relative weight ω on (a) UCMD and (b) CUB datasets respectively.

that the performance is relatively stable on UCMD and CUB in refs. [45, 85].

3) Effect of relative weight ω. In Figure 11, we show the effect of relative weight ω of L_{OSAP} on UCMD and CUB respectively. The relative weight ω in Equation (16) controls the weight between our two training losses L_{AP}^{SRF} and L_{OSP}: ω = 0 reduces L_{OSAP} to L_{AP}^{SRF} while ω = 1 to L_{OSP}. We can see in Figure 11 that training with the complete L_{OSAP} with both L_{AP}^{SRF} and L_{OSAP} is always better than using only one of the two losses. On UCMD and CUB datasets, the results both increase by ~2pt in the [0, 0.2] range and then decrease by ~2pt in the [0.2, 1] range. Consequently, the value of 0.2 is the optimal relative weight, achieving the best results on the two datasets.

5.5 Comparison with state-of-the-arts

In this section, we compare our OSAP-Loss to the recent AP approximation methods and some relevant state-of-the-art deep metric methods on three RSIR datasets and three CIR datasets. Firstly, we perform a comparison and analysis of AP approximation methods to verify the contribution in our paper. Secondly, we then conduct a further state-of-the-art comparison to demonstrate the superiority and effectiveness of our OSAP-Loss.

 ${\bf T}\,{\bf A}\,{\bf B}\,{\bf L}\,{\bf E}\,\,{\bf 7}$ Comparison between OSAP and state-of-the-art AP approximation methods

	CUB		SOP		INatu	ıralist	
Method	R@1	mAP@R	R@1	mAP@R	R@1	mAP@R	
FastAP [17]	58.9	22.9	78.2	51.3	53.5	19.6	
SoftBin [20]	61.2	24.0	80.1	53.5	56.6	20.1	
BlackBox [24]	62.6	23.9	80.0	53.1	52.3	15.2	
SmoothAP [25]	62.1	23.9	80.9	54.6	59.8	20.7	
ROADMAP [85]	64.2	25.3	82.0	56.5	64.5	25.1	
OSAP (ours)	67.1	26.5	78.7	51.4	66.3	25.4	
		UCMD	1	NWPUD	I	atternNet	
Method		mAP	r	nAP	r	nAP	
FastAP [17]		96.60	94.16		99.95		
SoftBin [20]	[20] 96.58		9	6.66	99.84		
BlackBox [24]		52.15	1	15.33		6.93	
SmoothAP [25]		96.49	9	2.53	9	9.93	
ROADMAP [85]		95.76	9	0.35	9	9.84	
OSAP (ours)		96.29	9	2.60	9	9.98	

Note: We report the results on six datasets. The results in **bold** indicate the best performance.

Abbreviations: AP, average precision; mAP, mean average precision; OSAP, Optimising samples after positive ones & average precision loss; SOP, Stanford Online Product.

					1		,		1										
		UCMI	0					NWPL	D					Patterr	Net				
Method	Venue	mAP	P@5	P@10	P@50	P@100	P@1000	mAP	$P(\underline{a})$	$P(\underline{a})$	P@50	P@100	P@1000	mAP	P@5	$P(\underline{a})$	P@50	P@100	P@1000
Triplet [75]	ICCV'17	92.96	98.04	96.63	92.62	46.16	4.69	93.82	98.65	96.85	96.07	94.83	15.34	94.94	99.52	97.92	96.13	95.07	15.61
N-pair-mc [69]	NIPS'16	91.81	94.04	91.46	90.49	45.08	4.67	93.06	97.86	95.12	94.35	98.15	15.46	94.11	97.94	95.15	94.33	98.17	15.52
ProxyNCA [70]	ICCV'17	95.72	97.98	96.65	94.23	47.02	4.71	97.68	97.54	97.57	97.91	97.44	15.69	97.71	98.56	98.69	98.89	98.45	15.74
LiftedStruct [43]	CVPR'16	96.08	98.90	97.82	95.78	47.46	4.76	97.47	97.03	97.42	97.63	97.72	15.76	98.58	98.05	98.62	98.75	98.88	15.79
[12] TISO	Entropy'19	97.34	98.98	98.42	96.93	48.67	4.86	98.54	99.05	98.15	96.34	98.45	15.78	98.52	99.09	98.03	96.68	98.69	15.83
DCL [2]	Remote Sensing'20	98.76	100.00	100.00	99.33	49.82	5.21	99.44	100.00	100.00	99.91	99.70	16.42	99.43	100.00	100.00	99.89	99.66	16.38
RRL [3]	TGRS'20	99.52	100.00	100.00	99.80	I	I	99.61	I	I	I	I	I	I	I	I	I	I	I
FastAP [17]	CVPR'19	96.61	100.00	100.00	99.99	78.99	7.90	94.38	95.74	95.54	95.12	94.52	13.88	99.95	99.96	99.96	99.96	90.96	55.89
BlackBox [24]	CVPR'20	52.70	86.40	80.52	58.98	44.25	7.80	15.35	40.65	35.25	22.53	17.45	6.24	4.17	2.63	2.63	2.63	2.63	2.63
SmoothAP [25]	ECCV'20	96.49	99.98	96.66	99.91	78.94	7.90	92.67	95.53	95.25	94.52	93.12	13.89	99.93	96.96	99.95	99.95	99.95	55.89
SoftBin* [20]	ICCV'19	96.58	100.00	100.00	99.98	78.99	7.90	94.05	95.85	95.66	95.20	94.30	13.69	99.95	79.97	79.97	70.00	79.97	55.92
ROADMAP [85]	NIPS'21	95.94	99.73	99.71	99.47	78.51	7.90	90.37	94.50	94.14	93.00	91.52	13.80	99.91	96.96	96.96	99.95	99.94	55.88
OSAP (w/o MFM)	This work	96.47	99.88	99.89	99.86	78.94	7.90	93.55	95.34	95.14	94.49	93.56	13.90	99.98	99.98	99.98	99.98	99.98	55.90
OSAP (w/MFM)	This work	99.64	100.00	100.00	99.94	79.54	7.94	99.76	100.00	100.00	99.95	99.84	13.97	99.99	10.00	100.00	100.00	100.00	56.31
Note: The embedding si: Abbreviations: mAP, me:	ze is 2048. We report the an average precision; MF	t results of M, memor	f mAP and ry-free mec	P@k. All 1 hanism; Oʻ	methods re SAP, Optir	ly on a star nising samp	dard convol les after pos	utional ba sitive ones	ckbone (ge & average	enerally Res precision	Net-50). T loss.	he results i	n bold indica	te the bes	t performa	nce.			

TABLE8 Comparison with state-of-the-art loss functions on UCMID [40], NWPUD [41] and PatternNet [1]

YUAN ET AL.

5.5.1 | Comparison to AP approximation methods

Here, we train all models using the same experimental setting (Backbone: ResNet-50; Batch size: 64; Embedding size: 512) as in Section 5.4. And, to compare our OSAP with recently proposed AP approximation methods, we conduct the experiments on three RSIR datasets and three CIR datasets. We report the results in Table 7. Specifically, we compare OSAP on three RSIR datasets and to recent AP approximation losses, including the soft-binning approaches FastAP [17] and SoftBin [20], BlackBox [24], SmoothAP [25], PNP-D_q [84] and ROADMAP [85]. We can see that OSAP outperforms most of the current AP approximation methods by a significant margin. The performance improvement is especially significant on large-scale datasets, such as the PatternNet dataset in RSIR and the INaturalist dataset in CIR.

5.5.2 | Comparison with state-of-the-arts

We use the same setting as in Section 5.4. For RSIR datasets, we only change the embedding size to 2048. For CIR datasets, we follow standard practices for ResNet-50 by using larger images (256×256 on SOP and CUB) and using max instead of average pooling and layer normalisation for CUB. We compare OSAP-Loss to other state-of-the-art methods across three RSIR datasets and three CIR datasets, and we report the results in Tables 8 and 9 respectively. We divide the comparison methods into two main categories: deep metric learning and AP approximation methods. OSAP-Loss falls in the second category. We show the results of OSAP-Loss with MFM and without MFM separately.

In Table 8, OSAP-Loss (w/o MFM) outperforms all previous methods on PatternNet dataset. We can see that the performances of the methods on the RSIR dataset are close to saturation. In Table 9, OSAP-Loss (w/o MFM) outperforms all previous AP approximation methods on CUB dataset and INaturalist dataset. It can be observed that MFM can increase the performance of OSAP-Loss by a large margin on all datasets, which further demonstrates the effectiveness of MFM.

5.5.3 | Visual assessment results

As a qualitative assessment, we show some retrieval results of OSAP-Loss on RSIR datasets in Figure 12. We show the query

TABLE 9 Comparison of state-of-the-art loss functions on SOP [43], CUB [42], and INaturalist [44]

				SOP			CUB				INati	uralist		
	Method	Venue	Dim	1	10	100	1	2	4	8	1	4	16	32
Deep metric learning	LiftedStruct [43]	CVPR'16	512	62.1	79.8	91.3	47.2	58.9	70.2	80.2	_	_	_	-
	Margin [75]	ICCV'17	512	72.7	86.2	93.8	63.6	74.4	83.1	90.0	58.1	75.5	86.8	90.7
	ProxyNCA [70]	ICCV'17	512	73.7	_	_	49.2	61.9	67.9	72.4	61.6	77.4	87.0	90.6
	MIC [53]	ICCV'19	512	77.2	89.4	95.6	66.1	76.8	85.6	-	-	-	_	-
	MS [54]	CVPR'19	512	78.2	90.5	96.0	65.7	77.0	86.3	91.2	-	_	_	_
	SoftTriple [76]	ICCV'19	512	78.3	90.3	95.9	65.4	76.4	84.5	90.4	-	_	_	_
	Circle [55]	CVPR'20	512	78.3	90.5	96.1	66.7	77.4	86.2	91.2	-	_	_	_
	SEC [56]	NIPS'20	512	78.7	90.8	96.6	68.8	79.4	87.2	92.5	_	_	_	_
	HORDE [57]	ICCV'19	512	80.1	91.3	96.2	66.8	77.4	85.1	91.0	-	_	_	_
	XBM [91]	CVPR'20	128	80.6	91.6	96.2	65.8	75.9	84.0	89.9	-	_	_	_
	Triplet SCT [58]	ECCV'20	512/64	81.9	92.6	96.8	57.7	69.8	79.6	87.0	-	_	_	_
AP-based loss	FastAP [17]	CVPR'19	512	76.4	89.0	95.1	_	_	_	_	60.6	77.0	87.2	90.6
	BlackBox [24]	CVPR'20	512	78.6	90.5	96.0	64.0	75.3	84.1	90.6	62.9	79.4	88.7	91.7
	SmoothAP [25]	ECCV'20	512	80.1	91.5	96.6	_	_	_	_	67.2	81.8	90.3	93.1
	SoftBin* [20]	ICCV'19	512	80.6	91.3	96.1	61.2	73.14	83.0	89.5	64.2	77.1	82.7	91.7
	$PNP-D_q$ [84]	AAAI'22	512	81.1	92.7	96.3	_	_	_	_	66.6	81.1	89.7	92.6
	ROADMAP [85]	NIPS'21	512	83.1	92.2	96.8	68.5	78.7	86.6	91.9	69.1	83.1	91.3	93.9
	OSAP (w/o MFM)	This work	512	79.9	91.3	96.5	69.8	79.6	88.0	92.8	70.7	84.3	92.5	94.7
	OSAP (w/MFM)	This work	512	84.4	93.1	97.3	70.5	80.2	88.3	93.2	71.0	84.7	92.8	94.9

Note: The embedding size is 512. We report the results of Recall@k. All methods rely on a standard convolutional backbone (generally ResNet-50). The results in bold indicate the best performance.



FIGURE 12 Illustration of image retrieval examples on three RSIR datasets. Golf course and baseball diamond in UCMD dataset (retrieval accuracy-99.64%), beach and desert in NWPUD (retrieval accuracy-99.76%), tennis court and football field in PtternNet dataset (retrieval accuracy-99.99%).

images, the most five similar retrieved images, and the corresponding retrieval accuracy. We can find that the semantic quality of the retrieval results is very good. These results also further demonstrate the effectiveness of our proposed OSAP-Loss.

6 | CONCLUSION

In this paper, we introduce a novel OSAP-Loss for RSIR, which improves the retrieval performance by directly optimising AP. Specifically, OSAP-Loss consists of three components: \mathcal{L}_{AP}^{SRF} \mathcal{L}_{OSP} , and MFM. To solve the weakness of sigmoid-based AP approximation, we proposed a SRF to replace the sigmoidbased ranking function, leading to the \mathcal{L}_{AP}^{SRF} being accurate. Afterwards, to overcome the non-decomposability in AP optimisation, we equip the \mathcal{L}_{AP}^{SRF} with \mathcal{L}_{OSP} to reduce the decomposability gap. Moreover, we develop an MFM to further thoroughly address the non-decomposability in AP optimisation, which sidesteps the constraint of the GPU memory to use large batch size training on a single GPU. We provide theoretical analysis as well as experimental results to demonstrate the superiority and effectiveness of OSAP-Loss. Extensive experiments that OSAP-Loss show the superiority and competitive performance on three RSIR datasets and three CIR datasets compared to the state-of-the-arts.

Furthermore, our OSAP-Loss is a data-driven method, which has the potential of propagating dataset biases. On some datasets, we find that the performance cannot be remarkable. We find that the performance cannot be exceptionally good on some datasets. OSAP-Loss is on par with or superior to existing loss functions. Future work plans to improve the robustness of our method and reduce the impact of bias in the dataset. We also plan to improve the stability of our approach and enhance its extensibility and applicability.

ACKNOWLEDGEMENT

This work was supported by the National Nature Science Foundation of China (No. U1803262, 62176191, 62171325), and Nature Science Foundation of Hubei Province (2022CFB018). This research was financially supported by fund from Hubei Province Key Laboratory of Intelligent Information Processing and Real-time Industrial System (Wuhan University of Science and Technology) (ZNXX2022001).

CONFLICT OF INTEREST

The author declares that there is no conflict of interest that could be perceived as prejudicing the impartiality of the research reported.

DATA AVAILABILITY STATEMENT

The data that support the findings of this study are openly available in UCMD at https://doi.org/10.1145/1869790. 1869829, reference number [39]. The data that support the findings of this study are openly available in NWPUD at https://doi.org/10.1109/JPROC.2017.2675998, reference number [40]. The data that support the findings of this study are openly available in PatternNet at https://doi.org/10.1016/j.isprsjprs. 2018.01.004, reference number [1]. The data that support the findings of this study are openly available in CUB at https://authors.library.caltech.edu/27452/1/CUB_200_2011.pdf, reference number [41]. The data that support the findings of this study are openly available in SOP at https://doi.org/10.1109/CVPR. 2016.434, reference number [42]. The data that support the findings of this study are openly available in INaturalist at https://doi.org/10.1109/CVPR.2018.00914, reference number [43].

ORCID

Xin Yuan https://orcid.org/0000-0003-3140-3243 *Xin Xu* https://orcid.org/0000-0003-0748-3669 *Xiao Wang* https://orcid.org/0000-0003-0770-9891 *Kai Zhang* https://orcid.org/0000-0003-0318-3255 *Liang Liao* https://orcid.org/0000-0002-2238-2420 *Zheng Wang* https://orcid.org/0000-0003-3846-9157 *Chia-Wen Lin* https://orcid.org/0000-0002-9097-2318

REFERENCES

- Zhou, W., et al.: Patternnet: a benchmark dataset for performance evaluation of remote sensing image retrieval. ISPRS J. Photogrammetry Remote Sens. 145, 197–209 (2018)
- Fan, L., Zhao, H., Zhao, H.: Distribution consistency loss for large-scale remote sensing image retrieval. Rem. Sens. 12(1), 175 (2020)
- Fan, L., Zhao, H., Zhao, H.: Global optimization: combining local loss with result ranking loss in remote sensing image retrieval. IEEE Trans. Geosci. Rem. Sens. 59(8), 7011–7026 (2020)
- Kang, J., et al.: Robust normalized softmax loss for deep metric learningbased characterization of remote sensing images with label noise. IEEE Trans. Geosci. Rem. Sens. 59(10), 8798–8811 (2021)
- Ge, W.: Deep metric learning with hierarchical triplet loss. In: Proceedings of the European Conference on Computer Vision, pp. 269–285 (2018)
- Chakrabarti, S., et al.: Structured learning for non-smooth ranking losses. In: Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 88–96 (2008)

24682222, 0, Downloaded from https://ietresearch.onlinelibrary.wiley.com/doi/10.1049/cit2.12151, Wiley Online Library on [2505/2023]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

- Caruana, R., et al.: Ensemble selection from libraries of models. In: Proceedings of the International Conference on Machine Learning, pp. 1–8 (2004)
- Yue, Y., et al.: A support vector method for optimizing average precision. In: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 271–278 (2007)
- Constantin, M.G., et al.: The predicting media memorability task at mediaeval 2019. In: Working Notes Proceedings of the MediaEval 2019 Workshop, pp. 27–29 (2019)
- Han, L., et al.: Cohesion intensive deep hashing for remote sensing image retrieval. Rem. Sens. 12(1), 101 (2019)
- Oksuz, K., et al.: Rank & sort loss for object detection and instance segmentation. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 3009–3018 (2021)
- Lin, Y., Lee, Y., Wahba, G.: Support vector machines for classification in nonstandard situations. Mach. Learn. 46(1), 191–202 (2002)
- Morik, K., Brockhausen, P., Joachims, T.: Combining statistical learning with a knowledge-based approach: a case study in intensive care monitoring. In: Proceedings of the International Conference on Machine Learning, pp. 268–277 (1999)
- Bartell, B.T., Cottrell, G.W., Belew, R.K.: Automatic combination of multiple ranked retrieval systems. In: Proceedings of the 17th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, pp. 173–181 (1994)
- Herschtal, A., Raskutti, B.: Optimising area under the roc curve using gradient descent. In: Proceedings of the International Conference on Machine Learning, pp. 49 (2004)
- Bruch, S., et al.: An analysis of the softmax cross entropy loss for learning-to-rank with binary relevance. In: Proceedings of the 2019 ACM SIGIR International Conference on Theory of Information Retrieval, pp. 75–78 (2019)
- Cakir, F., et al.: Deep metric learning to rank. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1861–1870 (2019)
- He, K., et al.: Hashing as tie-aware learning to rank. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4023–4032 (2018)
- He, K., Lu, Y., Sclaroff, S.: Local descriptors optimized for average precision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 596–605 (2018)
- Revaud, J., et al.: Learning with average precision: training image retrieval with a listwise loss. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 5107–5116 (2019)
- Ustinova, E., Lempitsky, V.: Learning deep embeddings with histogram loss. In: Proceedings of the Advances in Neural Information Processing Systems, pp. 4177–4185 (2016)
- Engilberge, M., et al.: Sodeep: a sorting deep net to learn ranking loss surrogates. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 10792–10801 (2019)
- Pogančić, M.V., et al.: Differentiation of blackbox combinatorial solvers. In: Proceedings of the International Conference on Learning Representations, pp. 1–19 (2020)
- Rolínek, M., et al.: Optimizing rank-based metrics with blackbox differentiation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7620–7630 (2020)
- Brown, A., et al.: Smooth-ap: smoothing the path towards large-scale image retrieval. In: Proceedings of the European Conference on Computer Vision, pp. 677–694 (2020)
- Liu, Z., et al.: Image retrieval on real-life images with pre-trained visionand-language models. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2125–2134 (2021)
- Sain, A., et al.: Stylemeup: towards style-agnostic sketch-based image retrieval. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8504–8513 (2021)
- Li, X., et al.: Qair: practical query-efficient black-box attacks for image retrieval. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3330–3339 (2021)

- Tadepalli, Y., et al.: Content-based image retrieval using Gaussian– Hermite moments and firefly and grey wolf optimization. CAAI Trans. Intell. Technol. 6(2), 135–146 (2021)
- Ahmad, F.: Deep image retrieval using artificial neural network interpolation and indexing based on similarity measurement. CAAI Trans. Intell. Technol. 7(2), 200–218 (2022)
- Deng, J., et al.: Masked face recognition challenge: the insightface track report. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 1437–1444 (2021)
- Meng, Q., et al.: Magface: a universal representation for face recognition and quality assessment. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 14225–14234 (2021)
- Zhu, Z., et al.: Webface260m: a benchmark unveiling the power of millionscale deep face recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 10492–10502 (2021)
- Liao, L., et al.: Exploiting effects of parts in fine-grained categorization of vehicles. In: Proceedings of the IEEE International Conference on Image Processing, pp. 745–749 (2015)
- Wang, Z., et al.: Re-identification = retrieval+ verification: back to essence and forward with a new metric. arXiv preprint, arXiv:201111506 (2020)
- Xie, P., et al.: Unsupervised video person re-identification via noise and hard frame aware clustering. In: 2021 IEEE International Conference on Multimedia and Expo, pp. 1–6 (2021)
- Jiang, M., et al.: Robust vehicle re-identification via rigid structure prior. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4026–4033 (2021)
- Xu, X., et al.: Rank-in-rank loss for person re-identification. ACM Trans. Multimed. Comput. Commun. Appl. 18, 1–21 (2022)
- Xu, X., et al.: Towards generalizable person re-identification with a bistream generative model. Pattern Recogn. 132, 108954 (2022)
- Yang, Y., Newsam, S.: Bag-of-visual-words and spatial extensions for land-use classification. In: Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems, pp. 270–279 (2010)
- Cheng, G., Han, J., Lu, X.: Remote sensing image scene classification: benchmark and state of the art. Proc. IEEE 105(10), 1865–1883 (2017)
- 42. Wah, C., et al.: The caltech-ucsd birds-200-2011 dataset (2011)
- OhSong, H., et al.: Deep metric learning via lifted structured feature embedding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4004–4012 (2016)
- Van Horn, G., et al.: The inaturalist species classification and detection dataset. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8769–8778 (2018)
- Cao, Z., et al.: Hashnet: deep learning to hash by continuation. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 5608–5617 (2017)
- Zhang, Z., et al.: Probability ordinal-preserving semantic hashing for largescale image retrieval. ACM Trans. Knowl. Discov. Data 15(3), 1–22 (2021)
- 47. Weng, Z., Zhu, Y.: Online hashing with bit selection for image retrieval. IEEE Trans. Multimed. 23, 1868–1881 (2021)
- Shen, F., et al.: Unsupervised deep hashing with similarity-adaptive and discrete optimization. IEEE Trans. Pattern Anal. Mach. Intell. 40(12), 3034–3044 (2018)
- Luo, Y., et al.: Robust discrete code modeling for supervised hashing. Pattern Recogn. 75, 128–135 (2018)
- Perronnin, F. et al.: Large-scale image retrieval with compressed Fisher vectors. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3384–3391 (2010)
- Arandjelovic, R., Zisserman, A.: All about VLAD. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1578–1585 (2013)
- Husain, S.S., Ong, E.J., Bober, M.: Actnet: end-to-end learning of feature activations and multi-stream aggregation for effective instance image retrieval. Int. J. Comput. Vis. 129(5), 1432–1450 (2021)
- Roth, K., Brattoli, B., Ommer, B.: Mic: mining interclass characteristics for improved metric learning. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 8000–8009 (2019)

- Wang, X., et al.: Multi-similarity loss with general pair weighting for deep metric learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5022–5030 (2019)
- Sun, Y., et al.: Circle loss: a unified perspective of pair similarity optimization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6398–6407 (2020)
- Zhang, D., Li, Y., Zhang, Z.: Deep metric learning with spherical embedding. In: Proceedings of the Advances in Neural Information Processing Systems, pp. 18772–18783 (2020)
- Jacob, P., et al.: Metric learning with horde: high-order regularizer for deep embeddings. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 6539–6548 (2019)
- Xuan, H., et al.: Hard negative examples are hard, but useful. In: Proceedings of the European Conference on Computer Vision, pp. 126–142 (2020)
- Zhao, H., et al.: Global-aware ranking deep metric learning for remote sensing image retrieval. Geosci. Rem. Sens. Lett. IEEE 19, 1–5 (2021)
- Kang, J., et al.: Rotation-invariant deep embedding for remote sensing images. IEEE Trans. Geosci. Rem. Sens. 60, 1–13 (2021)
- 61. Gu, M., et al.: Polsar ship detection based on a sift-like polsar keypoint detector. Rem. Sens. 14(12), 2900 (2022)
- 62. Mikriukov, G., Ravanbakhsh, M., Demir, B.: Deep unsupervised contrastive hashing for large-scale cross-modal text-image retrieval in remote sensing. *arXiv preprint, arXiv:220108125* (2022)
- Aptoula, E.: Remote sensing image retrieval with global morphological texture descriptors. IEEE Trans. Geosci. Rem. Sens. 52(5), 3023–3034 (2014)
- Scott, G.J., et al.: Entropy-balanced bitmap tree for shape-based object retrieval from large-scale satellite imagery databases. IEEE Trans. Geosci. Rem. Sens. 49(5), 1603–1616 (2011)
- 65. Hua, Y., et al.: Semantic segmentation of remote sensing images with sparse annotations. Geosci. Rem. Sens. Lett. IEEE 19, 1–5 (2021)
- Zhou, W., et al.: Learning low dimensional convolutional neural networks for high-resolution remote sensing image retrieval. Rem. Sens. 9(5), 489 (2017)
- Radenović, F., Tolias, G., Chum, O.: CNN image retrieval learns from bow: unsupervised fine-tuning with hard examples. In: Proceedings of the European Conference on Computer Vision, pp. 3–20 (2016)
- Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: a unified embedding for face recognition and clustering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 815–823 (2015)
- Sohn, K.: Improved deep metric learning with multi-class n-pair loss objective. In: Proceedings of the Advances in Neural Information Processing Systems, pp. 1857–1865 (2016)
- Movshovitz-Attias, Y., et al.: No fuss distance metric learning using proxies. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 360–368 (2017)
- Fan, L., et al.: Distribution structure learning loss (DSLL) based on deep metric learning for image retrieval. Entropy 21(11), 1121 (2019)
- Koestinger, M., et al.: Large scale metric learning from equivalence constraints. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2288–2295 (2012)
- Swets, D.L., Weng, J.J.: Using discriminant eigenfeatures for image retrieval. IEEE Trans. Pattern Anal. Mach. Intell. 18(8), 831–836 (1996)
- Harwood, B., et al.: Smart mining for deep metric learning. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2821–2829 (2017)
- Wu, C.Y., et al.: Sampling matters in deep embedding learning. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2840–2848 (2017)
- Qian, Q., et al.: Softtriple loss: deep metric learning without triplet sampling. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 6450–6458 (2019)
- Teh, E.W., DeVries, T., Taylor, G.W.: Proxynca++: revisiting and revitalizing proxy neighborhood component analysis. In: Proceedings of the European Conference on Computer Vision, pp. 448–464 (2020)
- Wang, X., et al.: Ranked list loss for deep metric learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5207–5216 (2019)

- Musgrave, K., Belongie, S., Lim, S.N.: A metric learning reality check. In: Proceedings of the European Conference on Computer Vision, pp. 681–699 (2020)
- Kim, S., et al.: Proxy anchor loss for deep metric learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3238–3247 (2020)
- Durand, T., Thome, N., Cord, M.: Exploiting negative evidence for deep latent structured models. IEEE Trans. Pattern Anal. Mach. Intell. 41(2), 337–351 (2018)
- McFee, B., Lanckriet, G.R.: Metric learning to rank. In: Proceedings of the International Conference on Machine Learning, pp. 1–8 (2010)
- Mohapatra, P., et al.: Efficient optimization for rank-based loss functions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3693–3701 (2018)
- Li, Z., et al.: Rethinking the optimization of average precision: only penalizing negative instances before positive ones is enough. In: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 1518–1526 (2022)
- Ramzi, E., et al.: Robust and decomposable average precision for image retrieval. In: Proceedings of the Advances in Neural Information Processing Systems, pp. 23569–23581 (2021)
- Law, M.T., Thome, N., Cord, M.: Learning a distance metric from relative comparisons between quadruplets of images. Int. J. Comput. Vis. 121(1), 65–94 (2017)
- Mohapatra, P., Jawahar, C., Kumar, M.P.: Efficient optimization for average precision SVM. In: Proceedings of the Advances in Neural Information Processing Systems, pp. 2312–2320 (2014)
- Suh, Y., et al.: Stochastic class-based hard example mining for deep metric learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7251–7259 (2019)
- Carvalho, M., et al.: Cross-modal retrieval in the cooking context: learning semantic text-image embeddings. In: The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, pp. 35–44 (2018)
- Faghri, F., et al.: Vse++: improving visual-semantic embeddings with hard negatives. arXiv preprint, arXiv:170705612 (2017)
- Wang, X., et al.: Cross-batch memory for embedding learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6388–6397 (2020)
- He, K., et al.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)
- Chopra, S., Hadsell, R., LeCun, Y.: Learning a similarity metric discriminatively, with application to face verification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 539–546 (2005)
- Hadsell, R., Chopra, S., LeCun, Y.: Dimensionality reduction by learning an invariant mapping. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1735–1742 (2006)
- Ye, F, et al.: Remote sensing image retrieval using convolutional neural network features and weighted distance. Geosci. Rem. Sens. Lett. IEEE 15(10), 1535–1539 (2018)
- Deng, J., et al.: Imagenet: a large-scale hierarchical image database. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255 (2009)
- Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization. arXiv preprint, arXiv:14126980 (2014)

How to cite this article: Yuan, X., et al.: OSAP-Loss: Efficient optimization of average precision via involving samples after positive ones towards remote sensing image retrieval. CAAI Trans. Intell. Technol. 1–22 (2023). https://doi.org/10.1049/cit2.12151