

Deep Angular Embedding and Feature Correlation Attention for Breast MRI Cancer Analysis

Luyang Luo¹, Hao Chen², Xi Wang¹, Qi Dou³,
Huangjing Lin¹, Juan Zhou⁴, Gongjie Li⁵, and Pheng-Ann Heng¹

¹ Dept. of Computer Science and Engineering,
The Chinese University of Hong Kong, Hong Kong SAR, China

² Insight Medical Technology, Co., Ltd., China

³ Dept. of Computing, Imperial College London, London, UK

⁴ Dept. of Radiology, The Fifth Medical Center
of Chinese PLA General Hospital, Beijing, China

⁵ Beijing Image Diagnostic Center of Rimag, Beijing, China

Abstract. Accurate and automatic analysis of breast MRI plays an important role in early diagnosis and successful treatment planning for breast cancer. Due to the heterogeneity nature, accurate diagnosis of tumors remains a challenging task. In this paper, we propose to identify breast tumor in MRI by Cosine Margin Sigmoid Loss (CMSL) with deep learning (DL) and localize possible cancer lesion by COReLation Attention Map (COAM) based on the learned features. The CMSL embeds tumor features onto a hyper-sphere and imposes a decision margin through cosine constraints. In this way, the DL model could learn more separable inter-class features and more compact intra-class features in the angular space. Furthermore, we utilize the correlations among feature vectors to generate attention maps that could accurately localize cancer candidates with only image-level label. We build the largest breast cancer dataset involving 10,290 DCE-MRI scan volumes for developing and evaluating the proposed methods. The model driven by CMSL achieved classification accuracy of 0.855 and AUC of 0.902 on the testing set, with sensitivity and specificity of 0.857 and 0.852, respectively, outperforming other competitive methods overall. In addition, the proposed COAM accomplished more accurate localization of the cancer center compared with other state-of-the-art weakly supervised localization method.

1 Introduction

Breast cancer is the most common malignancy affecting women worldwide [1]. Early diagnosis of breast cancer is essential for successful treatment planning, where Magnetic Resonance Imaging (MRI) plays a vital role for screening high-risk populations [2]. Clinically, radiologists use the Breast Imaging-Reporting and Data System (BI-RADS) to categorize breast lesions into different levels according to their phenotypic characteristics presented in MRI images, indicating

different degrees of cancer risk. However, such assessment suffers from inter-observer variance and often subjectively relies on the radiologists’ experience. Moreover, due to the heterogeneity nature, tumors of the same pathological result (malignant or benign) could have diverse patterns and hence result in different BI-RADS assessments. In other words, tumors could possess ambiguous inter-class difference and large intra-class variance, which poses a serious challenge to accurate diagnosis of breast cancer.

Generally, there are two major tasks regarding to breast MRI tumor analysis: identification of tumors and localization of cancer candidates. Recently, Deep Learning (DL) based approaches have demonstrated great potential in assisting diagnosis of breast cancer in an automatic and fast manner. Previous studies manually annotated tumors and deliberately extracted the corresponding slices or patches for classification [3,4]. Such methods depended on careful annotations both for training and testing and could not easily be adopted to clinical application. On the other hand, Guy et al. [5] proposed to first automatically localize the lesions and then classify cancer candidates at the second stage. Although the inference in the testing stage thereby was free of lesion delineation, these works still required annotations for model training. To get rid of manual extraction for region of interest (RoI), Gabriel et al. [6] proposed to meta-learn the breast MRI cancer classification problem with only image-level labels. However, all the mentioned studies were limited to small size datasets and consequently lack of generalization validation. More importantly, the relatively low precision or specificity reported in these works implied that the aforementioned problem of inter-class difference and intra-class variance has not been addressed yet.

To this end, we propose a Cosine-Margin Sigmoid Loss (CMSL) to tackle the heterogeneity problem for breast tumor classification and COrelation Attention Map (COAM) for precise cancer candidates localization, both with image-level labels only. The CMSL is extended from the cosine loss originally designed for face verification [7]. It embeds the deep feature vectors onto a hyper-sphere and learns a decision margin between classes in the angular feature space. As a result, the learned features possess more compact intra-class variance and more separable intra-class difference. In addition, we observe a RoI shifting problem of localizing cancer by class activation map [8]. Therefore, we propose a novel weakly supervised method, i.e., COAM, to localize cancer candidates more accurately by leveraging deep feature correlations based on the Gram matrix. Furthermore, we build the largest breast DCE-MRI dataset including 10,290 volume scans from 1715 subjects to develop and evaluate our methods.

2 Methods

Our framework of breast MRI tumor analysis consists of two parts as illustrated in Fig. 1. One is tumor classification by deep angular embedding driven DL

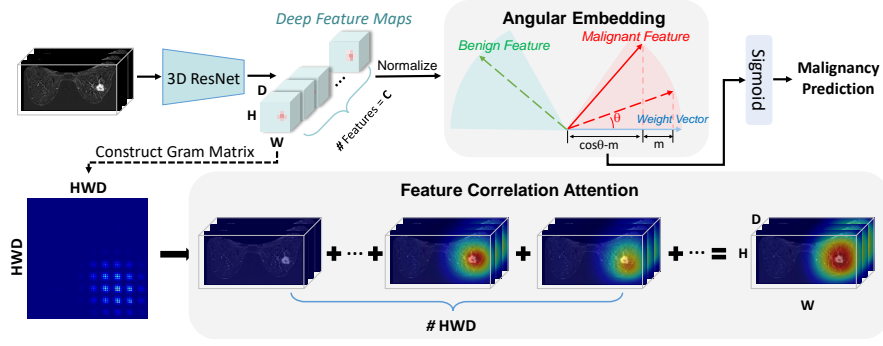


Fig. 1. The framework of breast MRI cancer analysis. A 3D ResNet is first trained with CMSL by embedding the deep features onto hyper-sphere. In the testing stage, the deep features are used to construct Gram matrix to obtain correlation attention map.

network. The other is weakly supervised cancer candidates localization with feature correlation attention map.

2.1 Cosine Margin Sigmoid Loss for Tumor Classification

The phenotype of tumors has ambiguous inter-class difference and large intra-class variance. Accordingly, the features learned by the DL model could inherit these characteristics. To address this issue, we start by revisiting the traditional sigmoid loss for binary classification problem. Given the input feature vector x of the last fully connected (FC) layer and its corresponding label y , the binary sigmoid loss is as follows:

$$\mathcal{L}(w; x) = -y \cdot \log(p(y | x)) - (1 - y) \cdot \log(1 - p(y | x)) \quad (1)$$

$$= -y \cdot \log\left(\frac{1}{1 + e^{-w^T x}}\right) - (1 - y) \cdot \log\left(1 - \frac{1}{1 + e^{-w^T x}}\right) \quad (2)$$

where w is the weight parameter of the FC layer, and $p(y | x)$ represents the probability of x being classified to y . To distinguish different classes, the DL model is expected to give different predictions by adjusting the value of $w^T x$. Notice that $w^T x = \|w\| \|x\| \cos\theta$, where θ is the angle between feature vector x and weight vector w , and $\|\cdot\|$ is the L_2 norm operation. Generally, the DL model would implicitly alter $\|w\|$ and $\|x\|$ in the Euclidean space and $\cos\theta$ in the angular space. However, the aforementioned heterogeneity issue could lead to ambiguous features that are quite hard to discriminate. To this end, constraints on feature distances are considered to regulate the DL model for more separable inter-class features and more compact intra-class features [7]. Since Euclidean distance is not bounded and hence difficult to constrain, we prefer to add regularization on the angular distance which is bounded by $-1 \leq \cos\theta \leq 1$. Specifically, we eliminate the influence of the norms $\|x\|$ and $\|w\|$ by modifying

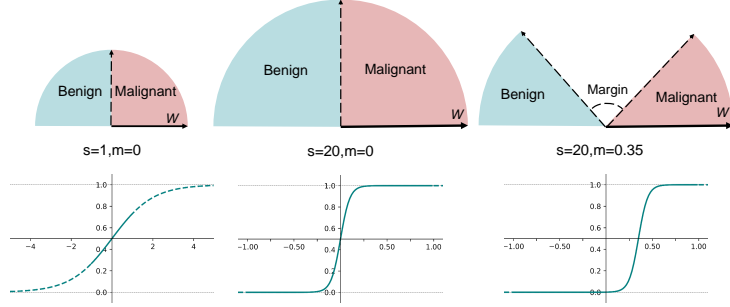


Fig. 2. The illustration of NSL with $s = 1$, NSL with $s = 20$ and CMSL with $s = 20$ and $m = 0.35$. First row is the geometric interpretation of feature projection on a 2D sphere. Dashed arrows represent the decision boundaries. Second row is the plot of corresponding sigmoid functions. Dashed curves represent the values out of range.

the computation of $p(y | x)$ to:

$$p(y | x) = \frac{1}{1 + e^{-s \frac{w^T x}{\|w\| \|x\|}}} = \frac{1}{1 + e^{-s \cdot \cos\theta}} \quad (3)$$

where s is a hyper-parameter adjusting the slope of the sigmoid function and controlling the back propagated gradient values. If s is too small, the loss cannot converge to 0 because the sigmoid function is not able to reach its saturation area, given that $-1 \leq \cos\theta \leq 1$. On the contrast, if s is set to a large value, the sigmoid function could easily reach the saturation area and result in small gradients, which prevents the network from learning sufficient knowledge. Following [7], we refer to the loss with modified p in Eq.(3) as Normalized Sigmoid Loss (NSL), which focuses on separating features in the angular space with decision boundary $\cos\theta = 0$ for both classes. Geometrically, we embed the feature vector and the weight vector onto a hyper-sphere whose radius is tuned by s . However, the ambiguous features can still distribute near this boundary. Therefore we add an explicit guidance to NSL as follows:

$$\mathcal{L}(w; x) = -y \cdot \log\left(\frac{1}{1 + e^{-s \cdot (\cos\theta - I(y) \cdot m)}}\right) - (1 - y) \cdot \log\left(1 - \frac{1}{1 + e^{-s \cdot (\cos\theta - I(y) \cdot m)}}\right) \quad (4)$$

where $I(\cdot)$ is an indicator function. $I(y) = 1$ if $y = 1$ and $I(y) = -1$ otherwise. m is a hyper-parameter that changes the decision boundaries for separating two classes (0 and 1 for benign and malignant) to: $B_0 : \cos\theta + m < 0$ and $B_1 : \cos\theta - m > 0$. Hence a decision margin is imposed by m in the angular space to make the learned inter-class features more separable. Consequently, the distribution space of features shrinks, which eventually leads to more compact intra-class features. Fig. 2 shows a comparison among different sigmoid functions and the corresponding geometric illustrations.

2.2 Feature Correlation Attention for Cancer Localization

Based on the well trained network, localization of cancer candidates can provide more evidences for clinical reference. Therefore, our secondary goal is to localize possible cancers out of other lesion mimics. It is natural for DL studies to use Class Activation Map (CAM) [8] for obtaining the Region of Interest (RoI) when only image-level label is available. However, it can not be well generalized to our case due to an observed RoI shifting problem. With the CNN going deeper, the reception fields of neurons become larger, hence neighbors of the tumor feature also capture views over the tumor patch in the image. Since the feature vectors corresponding to different classes could be ambiguous, the classifier layer would possibly tend to find discriminative patterns in the neighbors. Consequently, the corresponding RoI generated by CAM would shift from the desired target.

To tackle this problem, we further figure out two insights of our task. First, the feature vectors of the same semantic (malignant or normal) ought to have high correlations with each other. Second, through a series of rectified linear units, the network would implicitly learns large activation values for features related to suspicious cancer patch (with label “1”), and small activation values for features related to normal patch (with label “0”). Based on these two intuitions, we leverage the Gram matrix [9] to find the RoI. Given the deep feature map $X \in \mathbb{R}^{H \times W \times S \times C}$ from the last activation layer, where H, W, S and C are the height, width, number of slices and number of channels, respectively, we first reshape X to $X' \in \mathbb{R}^{N \times C}$, where $N = H \times W \times S$. Then we compute an attention vector $M \in \mathbb{R}^N$ as follows:

$$M_i = \sum_{j=1}^N G_{i,j} = \sum_{j=1}^N \sum_{k=1}^C X'_{i,k} X'_{j,k} \quad (5)$$

where $G \in \mathbb{R}^{N \times N}$ is the Gram matrix over the set of deep feature vectors in X' . Each entry $G_{i,j}$ is the inner product of X'_i and X'_j , representing the correlation between i -th and j -th vector. Because our network is trained for binary classification, it enables the gap between large and small activation values of feature vector related to suspicious cancer and normal patch. Correspondingly, the correlation value would also be relatively large or small according to the activation values of the features. Inspired by [10], each column G_i can be interpreted as a sub-attention map implying the network’s attention of the class that i -th vector belongs to. Thus the above operation is equal to element-wise summation over all sub-attention maps G_i . Moreover, since G is symmetric, the element-wise summation is also equivalent to summing over G_i to be the value of M_i . Essentially, $\sum_{j=1}^N G_{i,j}$ indicates the *importance* of i -th feature determined by the sub-attention of the feature map at its i -th position. At last, by simply reshape M to $H \times W \times S$ we are able to obtain an attention map purely based on the deep feature correlations. We refer to this method as CORelation Attention Map (COAM). It is worth mentioning that COAM is related to the self-attention mechanism [10] and the stationary feature space representation [9].

However, it differs from these works because the Gram matrix is not involved at any optimization stage and directly used for attention map generation.

3 Experiments and Results

3.1 Implementation Details

Dataset We built the largest breast tumor Dynamic Contrast Enhanced (DCE) MRI dataset involving 10,290 scans from 1715 subjects, with 1137 cases containing malignant tumors and 578 cases containing benign tumors. All of the scans were conducted with a 1.5-T Siemens system. We collected 6 DCE-MRI subtraction scans and 1 non-fat suppressed T1 scan from each subject. BI-RADS categories were assessed by 3 radiologists. Pathological labels was given by biopsy or surgery diagnosis. The data were randomly divided into training, validation and testing sets with 1204, 165 and 346 subjects, respectively.

Preprocessing Frangi’s approach[11] was first applied on the slices of each non-fat suppressed T1 scan to detect evident edges. Next, thresholding, small connected component removal and hole filling were sequentially employed to obtain coarse breast region masks. Afterwards, the 2D masks were stacked into volumes, followed by Gaussian smooth. We used the 3D masks to segment the subtractions. Note that the DCE-MRI and non-fat suppressed scans were originally registered in the scanning machine. Finally we clipped and normalized the intensity values, concatenated 6 subtractions, and cropped or padded the data to a fixed size of $340 \times 220 \times 128$ as the model inputs.

Training Strategy We used 3D ResNet34 [12] as the base model and replaced the global average pooling layer and FC layer with an $1 \times 1 \times 1$ convolutional layer appended with a pooling layer. The hyper-parameter s and m were set to 20 and 0.35, respectively, similar to [7]. The learning rate was initially set to 10^{-4} and decreased 10 times when training error stagnated. The base model is trained until convergence and then employed to initialize all other methods.

3.2 Evaluation and Comparison

Tumor Classification We conducted comparison among several deep learning methods: (1)*2D MIL*: a multi instance method aggregating features from 2D slices by 2D ResNet34 [13]; (2)*3D ResNet*: a 3D implementation of ResNet34;

Table 1. Comparison of different methods on cancer classification.

Method	Accuracy	Sensitivity	Specificity	F1	AUC
2D MIL [13]	0.789	0.870	0.626	0.846	0.842
3D ResNet [12]	0.821	0.840	0.783	0.862	0.880
3D Sparse MIL [14]	0.832	0.857	0.783	0.872	0.885
3D DK-MT [15]	0.824	0.896	0.643	0.864	0.883
3D ResNet+NSL	0.821	0.840	0.783	0.862	0.874
3D ResNet+CMSL (ours)	0.855	0.857	0.852	0.888	0.902

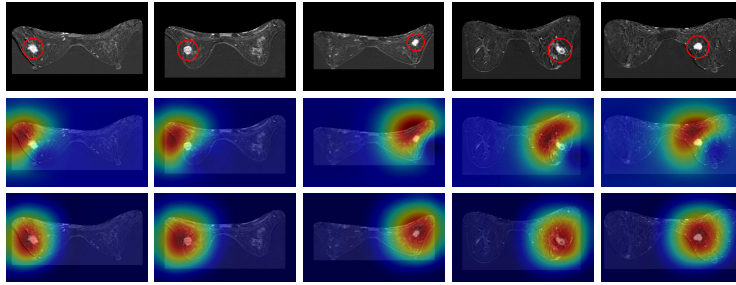


Fig. 3. Comparison Between CAM and COAM. We select typical slices from different subjects for a qualitative demonstration. First row: DCE-MRI subtraction slice; second row: visualization of CAM; third row: visualization of COAM. Cancer lesions are circles by red. Best viewed in color.

(3) *3D Sparse MIL*: a sparse label assign method [14]; (4) *3D DK-MT*: a domain knowledge driven multi-task learning network [15]; (5) *3D ResNet+NSL*: Normalized sigmoid loss based on (2); (6) *3D ResNet+CMSL*: our proposed CMSL based on (2). We computed the accuracy, specificity, sensitivity, F1 score and AUC as the evaluation metrics. Experimental results are reported in Table 1.

Compared with 2D method, 3D approaches achieved better results by utilizing more spatial information. Both *3D Sparse MIL* and *3D DK-MT* adopted additional assumption or knowledge, leading to better performance than vanilla *3D ResNet*. Noticeably, *3D DK-MT* showed poor specificity, which is possibly due to imbalanced auxiliary knowledge (more BI-RADS 4 and 5 than 3) that dominated the learning process. For deep angular embedding based methods like *3D ResNet+NSL*, simply taking the features into angular space without margin constraint caused certain performance decay. It implied that the network cannot learn sufficient knowledge if s is set to a large value. Moreover, our proposed *3D ResNet+CMSL* method significantly improved the results. The underlying reason is that it could learn more discriminative patterns by imposing cosine margin. Our method achieved the highest specificity with over 7.9% better than all other methods and kept a comparable sensitivity at the mean time. It exceeded all other methods with over 2% in AUC, over 3% in accuracy and over 1.5% in F1 score, proving that addressing the inter- and intra-class problem can improve performance of breast tumor classification.

Cancer Localization To evaluate the performance of COAM, we invited the radiologists to manually annotate 85 samples that were classified as malignant by our model. We compared our method with CAM by computing the Euclidean distance between center position of the annotation and the voxel position with highest value in the attention map. Then the distance is multiplied by the voxel spacing, i.e., 1.1 mm, as the final measurement. The criteria is reported in the form of $mean \pm std$, where $mean$ and $stdv$ represent the mean value and standard deviation of the center distances over 85 samples, respectively. Compared to the distance of 39.84 ± 8.82 mm by CAM, COAM showed a significant advantage with

18.26±13.65 mm only. Fig. 3 showed the qualitative comparison with these two methods.

4 Conclusion

In this paper, we propose the cosine margin sigmoid loss for breast tumor classification and correlation attention map for weakly supervised cancer candidates localization based on MRI scans. First, we use CMSL driven deep network to learn more separable inter-class features and more compact intra-class features which effectively tackle the heterogeneity problem of tumors. In addition, the proposed COAM leverages correlations among deep features to localize region of interests in a weakly supervised manner. Extensive experiments on our large-scale dataset demonstrates the efficacy of our methods which outperform other state-of-the-art approaches significantly on both tasks. Our methods are general and can be extended to many other fields.

References

1. DeSantis, C. E., et al.: Breast cancer statistics, 2017, racial disparity in mortality by state. In: *CA: a cancer journal for clinicians* **67**(6). pp439-448. (2017).
2. Kuhl, C., et al.: Prospective multicenter cohort study to refine management recommendations for women at elevated familial risk of breast cancer: the EVA trial. *J Clin Oncol* 28.9 (2010): 1450-1457.
3. Zheng, H., et al.: Small Lesion Classification in Dynamic Contrast Enhancement MRI for Breast Cancer Early Detection. *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, Cham. (2018).
4. Amit, G., et al.: Classification of breast MRI lesions using small-size training sets: comparison of deep learning approaches. In: *Medical Imaging 2017: Computer-Aided Diagnosis*. Vol. 10134. International Society for Optics and Photonics, 2017.
5. Amit, G., et al.: Hybrid mass detection in breast MRI combining unsupervised saliency analysis and deep learning. *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, Cham. pp. 594-602. (2017).
6. Maicas, G., et al.: Training medical image analysis systems like radiologists. *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, Cham. (2018).
7. Wang, H., et al.: Cosface: Large margin cosine loss for deep face recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2018).
8. Zhou, B., et al.: Learning deep features for scene recognition using places database. *Advances in neural information processing systems*. pp. 487-495. (2014).
9. Gatys, L., Ecker, A. S., & Bethge, M.: Texture synthesis using convolutional neural networks. *Advances in neural information processing systems* (pp. 262-270). (2015).
10. Fu, J., et al.: Dual attention network for scene segmentation. *arXiv preprint arXiv:1809.02983* (2018).
11. Frangi, A. F., et al.: Multiscale vessel enhancement filtering. *International conference on medical image computing and computer-assisted intervention*. pp. 130-13. Springer, Berlin, Heidelberg. (1998).

12. He, K., et al: Deep residual learning for image recognition. Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770-778. (2016).
13. Wu, J., et al.: Deep multiple instance learning for image classification and auto-annotation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2015).
14. Zhu, W., et al.: Deep multi-instance networks with sparse label assignment for whole mammogram classification. International Conference on Medical Image Computing and Computer-Assisted Intervention. Springer, Cham, pp. 603-611.(2017).
15. LIU, J., et al.: Integrate Domain Knowledge in Training CNN for Ultrasonography Breast Cancer Diagnosis. International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 868-875. Springer, Cham. (2018).