# EDGEMIXUP: Embarrassingly Simple Data Alteration to Improve Lyme Disease Lesion Segmentation and Diagnosis Fairness

Haolin Yuan<sup>§</sup>, John Aucott<sup>†</sup>, Armin Hadzic<sup>†</sup>, William Paul<sup>†</sup>, Marcia Villegas de Flores<sup>§</sup>, Philip Mathew<sup>§†</sup>, Philippe Burlina<sup>†</sup>, Yinzhi Cao<sup>§</sup>

<sup>§</sup>Johns Hopkins University

<sup>†</sup>Johns Hopkins Applied Physics Lab, <sup>‡</sup>Johns Hopkins University School of Medicine {hyuan4, mvilleg5, yinzhi.cao}@jhu.edu, {william.paul, philip.mathew, Philippe.Burlina}@jhuapl.edu armin.hadzic.scholar@gmail.com, jaucott2@jhmi.edu

Abstract. Lyme disease is a severe skin disease caused by tick bites, which affects hundreds of thousands of people. One task in diagnosing Lyme disease is lesion segmentation, i.e., separating benign skin from lesions, which can not only help clinicians to focus on lesions but also improve downstream tasks such as disease classification. However, it is challenging to segment Lyme disease lesions due to the lack of wellsegmented, labeled Lyme datasets and the nature of Lyme, e.g., the typical bull's eye lesion and its closeness to normal skin. In this paper, we design a simple yet novel data preprocessing and alteration method, called EDGEMIXUP, to help segment Lyme lesions on imbalanced training datasets. The key insight is to deploy a linear combination of lesion edge, either detected or computed, and the source image highlights the affected lesion area so that a learning model focuses more on the preserved lesion structure instead of skin tone, thus iteratively improving segmentation performance. Additionally, the improved edge from lesion segmentation can be further used for Lyme disease classification—e.g., in differentiating Lyme from other similar lesions including tinea corporis and herpes zoster-with improved model fairness on different subpopulations.

## 1 Introduction

Medical Image Analysis has greatly benefited from advances in AI [1] yet some improvements still remain to be addressed, importantly in areas that allow both algorithmic performance and fairness [2], and in certain medical applications that promise to significantly lessen morbidity and mortality. Early detection of skin lesions is such an endeavor as it can aid in identifying infectious diseases with cutaneous manifestations. Lyme disease is an example of that with a potentially diagnostic skin lesion [3]—which is caused by the bacterium *Borrelia burgdorferi* and leads to nearly 476,000 cases per annum during 2010–2018 [4]. The earliest and most treatable phase of Lyme disease is manifested via a red concentric lesion at the site of a tick bite, called erythema migrans (EM) [5]. While the

#### 2 Yuan et al.

EM pattern may appear simple to recognize, its diagnosis can be challenging for those with or without a medical background alike, as only 20% of United States patients have the stereotypical bull's eye lesion [6]. When skin lesions are atypical they can be mistaken for other diseases such as tinea corporis (TC) or herpes zoster (HZ), two other diseases acting as confusers for Lyme, considered herein. This has increased interest in medical applications of deep learning (DL), and using deep convolutional neural networks (CNNs), to assist clinicians in timely and accurate diagnosis of conditions including Lyme disease, TC and HZ [7–9].

One important diagnosis task is to segment Lyme lesion, particularly the EM pattern, from benign skins. Such DL-assisted segmentation not only helps clinicians in pre-screening patients but also improves downstream tasks such as lesion classification. However, while Lyme disease lesion segmentation is intuitively simple, it is challenging due to the following reasons. First, there lacks of a well-segmented dataset with manual labels on Lyme disease. On one hand, some datasets—such as HAM10000 [10] and ISBI Challenges [11]—have manual annotated segmentations for diseases like melanoma, but they do not have Lyme disease lesions. On the other hand, some datasets—such as Groh et al. [12]—have Lyme disease and skin tone and classification labels, but not segmentation.

Second, the segmentation of Lyme lesion is itself challenging due to the nature of EM pattern. Specifically, a typical Lyme lesion exhibits a bull's eye pattern with one central redness and one outer circle, which is different from darkness lesion in cancer-related skin disease like melanoma. Furthermore, clinical data collected for training is usually imbalanced in some properties, e.g., more samples with light skins compared with dark skins. Therefore, existing skin disease segmentation [13] as well as existing general segmentation works, such as U-Net [14], polar training [15], ViT-Adapter [16], and MFSNet [17], usually suffer from relatively low performance and reduced fairness [2, 18, 19].

In this paper, we present the first Lyme disease dataset that contains labeled segmentation and skin tones. Our Lyme disease dataset contains two parts: (i) a classification dataset, composed of more than 3,000 diseased skin images that are either obtained from public resources or clinicians with patient-informed consent, and (ii) a segmentation dataset containing 185 samples that are manually annotated for three regions—i.e., background, skin (light vs. dark), and lesion—conducted under clinician supervision and Institutional Review Boards (IRB) approval. Our dataset with manual labels is available at this URL [20].

Secondly, we design a simple yet novel data preprocessing and alternation method, called EDGEMIXUP, to improve Lyme disease segmentation and diagnosis fairness on samples with different skin-tones. The key insight is to alter a skin image with a linear combination of the source image and a detected lesion boundary so that the lesion structure is preserved while minimizing skin tone information. Such an improvement is an *iterative* process that gradually improves lesion edge detection and segmentation fairness until convergence. Then, the detected, converged edge in the first step also helps classification of Lyme diseases via mixup with improved fairness. Our source code is available at this URL [20].



(a) Original image (b) Heat map for EDGEMIXUP (c) Heat map for legacy analysis

We evaluate EDGEMIXUP for skin disease segmentation and classification tasks. Our results show that EDGEMIXUP is able to increase segmentation utility and improve fairness. We also show that the improved segmentation further improves classification fairness as well as joint fairness-utility metrics compared to existing debiasing methods, e.g., AD [21] and ST-Debias [22].

## 2 Motivation

In this section, we motivate the design of EDGEMIXUP by showing that added lesion boundary helps a DL model focus more on the lesion part instead of other features such as skin or background. Note that not all skin disease datasets are carefully processed either due to the large amount of work required or the scarcity of data samples collected, e.g., SD-198 [23] contains samples that are taken under variant environments. Specifically, we train two ResNet-34 models using the same dataset with and without EDGEMIXUP for a classification task of skin disease. We keep all hyper-parameters exactly the same for two models, and only augment the same image with and without mixing lesion boundary up with the original image. We generate initial lesion edges using EdgeMixup, which we will elaborate in following sections. Figure 1 shows the original image (Figure 1a) as well as two models' attention as heat-maps where red color represents the highest attention, yellow a higher attention, and purple the least attention. EDGEMIXUP helps the model to focus more on the lesion area comparing Figure 1b and 1c. The reason is that a legacy diagnosis has no information about lesion and does not know where to locate its focus, thus easily gets distracted by fingers instead of the lesion pattern.

# 3 Method

In this section, we first give the definition for model fairness, and we then describe the design of EDGEMIXUP for the purpose of de-biasing in Figure 2 and

Fig. 1: A motivating example to illustrate why EDGEMIXUP improves model performance and reduces biases via mixing up lesion boundary with original image (Heatmap is generated via Grad-CAM).



Fig. 2: EdgeMixup Process

Algorithm 1. We consider any model f, either a classification model  $f_{class}$  or a segmentation model  $f_{seg}$ , to be biased against certain skin-tone  $st_2$  if given metrics M and samples  $x_{st_1}$  and  $x_{st_2}$  from class y, where  $st_1$  and  $st_2$  are different skin-tones according to their ITA scores,  $M(f(x_{st_1}), y) > M(f(x_{st_2}), y)$ . If there exists a model f such that  $M(f(x_{st_1}), y) = M(f(x_{st_2}, y))$ , we consider it perfectly fair for  $st_1$  and  $st_2$  skin-tone samples.

EDGEMIXUP improves model fairness on light and dark skin samples in both segmentation and classification tasks, and it has two major components: (i) edge detection using mixup, and (ii) data preprocessing and alteration for downstream tasks. More specifically, our proposed edge detection has two parts: initial edge detection and iterative improvement.

Initial Edge Detection: The purpose of initial detection, which is documented in the Initial edge detection function of Algorithm 1, is to provide a starting point, i.e., a rough boundary, for the next step of iterative improvement. The high-level idea is that EDGEMIXUP detects several edge candidates using the color range of ground-truth lesions in both Red-Green-Blue (RGB) and Hue-Saturation-Value (HSV) color space and then selects the target edge using a learning model based on the output confidence score. First, EDGEMIXUP trains a classification model based on a mixup of the ground-truth segmentation under clinician supervision and the original image (Line 7). Second, EDGEMIXUP generates many edge candidates. For example, EDGEMIXUP collects the mean range of lesion color from the training set and use the range as threshold to filter out any given sample for a candidate mask (Line 9). Lastly, EDGEMIXUP selects an edge candidate with the highest confidence score output by the learning model (Line 11) and returns it as the edge for this given sample. Note that the initial edge detection is irrelevant to the sample size of a particular subpopulation, thus improving the fairness. That is, even if the original dataset is imbalanced, as long as one sample from a subpopulation exists, the color range of the sample's lesion is considered in the initial detection.

**Iterative Edge Improvement:** EDGEMIXUP includes iterative edge improvement in the training phase of our segmentation model to further improve model utility. The intuitive reason of utilizing such algorithm is that by applying the

#### Algorithm 1 Pseudo-code of EDGEMIXUP

**Input:** A labelled sample  $(x, y) \in D$ , mixup weights  $\alpha$ , ground-truth edged training set  $D_{edge_gt}^{train}$ **Output:** dataset  $D_{\texttt{final}_{edge}}$  in which each sample has it lesion edge highlighted  $(x_{\texttt{edge}}, y)$ 1: function main() 2:  $D_{\texttt{initial\_edge}} = \texttt{Initial\_edge\_detection}(D, \alpha)$ 3:  $D_{\texttt{final\_edge}} = \texttt{Iterative\_edge\_improvement}(D_{\texttt{initial\_edge}}, \alpha)$ 4: return D<sub>final edge</sub> 5: end function 6: function Initial edge detection $(D, \alpha)$ Train classification model  $m_{class}$  using  $D_{edge_gt}^{train}$ 7: 8: for each sample  $x \in D$  do 9: Get all edge candidates  $\{edge_1, edge_2, .., edge_n\}$  for each sample x 10: Mixup each edge candidate with x11: Query  $m_{\text{class}}$  using all mixed-up  $\{x_{\text{edge}_1}, \dots x_{\text{edge}_n}\}$  and choose the optimal edge  $\text{edge}_{\text{opt}}$ 12:Generate edged sample  $x_{edge} = \text{Mixup}(x, edge_{opt}, \alpha)$ 13:end for return  $D_{edge}$ 14:15: end function 16: function <code>Iterative\_edge\_improvement(D\_{edge}, \alpha)</code> Train the first model  $\overline{m_{\text{iter}}}$  using edged dataset  $D_{\text{edge}}^{\text{train}}$ 17:18: Evaluate  $m_{\text{iter}}$  using  $D_{\text{edge}}^{\text{test}}$  and get current\_Jaccard  $\begin{array}{l} \texttt{best}\_\texttt{Jaccard}=0\\ \texttt{iter}=1 \end{array}$ 19:20: 21: while current Jaccard > best Jaccard do 22: best Jaccard = current Jaccard 23:Predict lesion masks using  $m_{iter}$ , convert them to lesion edge edge24: Generate new training set for next model  $Mixup(D_{train}, edge, \alpha)$ Train a model for next iteration  $m_{\text{iter}+1}$ Evaluate  $m_{\text{iter}+1}$  using edged  $D_{\text{edge}}^{\text{test}}$  and get current\_Jaccard 25:26:27:iter += 128:end while 29: end function 30: function  $Mixup(x, edge, \alpha)$ 31: return  $(\alpha \cdot x + (1 - \alpha) \cdot edge)$ 32: end function



Fig. 3: Illustration of iterative edge improvement on different iterations with train loss

mixup of detected edge and original image, given the lesion boundary feature detected in the previous iteration, the next-iteration segmentation model can converge better and the lesion boundary predicted by it is fine-grained. Specifically, EDGEMIXUP iteratively trains segmentation models from scratch, and we let the model trained in the previous iteration to predict lesion edge, which is then mixed-up with original training samples as the new training set for next-iteration model. The high-level idea is that when the lesion is restricted in a small

Table 1: Annotated segmentation and classification dataset characteristics, broken down by ITA-based skin tones (light skin/ dark skin) and disease types.

Split	Skin				SD-sub						
	NO	EM	ΗZ	TC	Total	DF	KA	$\mathbf{PG}$	TC	TF	Total
seg		$\begin{array}{c} 62 \\ 47/15 \end{array}$	$\begin{array}{c} 62 \\ 46/16 \end{array}$	$\begin{array}{c} 61 \\ 40/21 \end{array}$	$\begin{array}{c} 185\\ 133/52 \end{array}$	$\begin{vmatrix} 30 \\ 23/7 \end{vmatrix}$	$\begin{array}{c} 30\\27/3\end{array}$	$\begin{array}{c} 30\\27/3\end{array}$	$\begin{array}{c} 30\\24/6\end{array}$	$\begin{array}{c} 30\\29/1\end{array}$	$\begin{array}{c}150\\130/20\end{array}$
class	$\begin{smallmatrix} 885\\822/63 \end{smallmatrix}$	$\begin{array}{c} 740 \\ 682/58 \end{array}$	$\begin{array}{c} 698 \\ 608/90 \end{array}$	$\begin{array}{c} 704 \\ 609/95 \end{array}$	$3027 \\ 2721/306$	$\begin{vmatrix} 40 \\ 36/4 \end{vmatrix}$	$\begin{array}{c} 40\\ 36/4\end{array}$	$\begin{array}{c} 40\\29/1\end{array}$	$\begin{array}{c} 40\\ 33/7\end{array}$	$\begin{array}{c} 40\\ 30/0 \end{array}$	$\begin{array}{c} 200\\ 164/16 \end{array}$

affected area, further detection will refine and constrain the detected boundary. Besides, EDGEMIXUP calculates a linear combination of original image and lesion boundary, i.e., by assigning the weight of original image as  $\alpha$  and lesion boundary as  $1 - \alpha$ . Figure 3 shows the edge-mixed-up images for different iterations. EDGEMIXUP removes more skin areas after each iteration and gradually gets close to the real lesion at the third iteration.

#### 4 Datasets

We present two datasets: (i) a dataset collected and annotated by us (called Skin), and (ii) a subset of SD-198 [23] with our annotation (called SD-sub). First, We collect and annotate a dataset with 3,027 images containing three types of disease/lesions, i.e., Tinea Corporis (TC), Herpes Zoster (HZ), and Erythema Migrans (EM). All skin images are either collected from publicly available sources or from clinicians with patient informed consent. Then, a medical technician and a clinician in our team manually annotate each image. For the segmentation task, we annotate skin images into three classes: background, skin, and lesion; then, for the classification task, we annotate skin images by classifying them into four classes: No Disease (NO), TC, HZ, and EM. We name it as Skin-class for later reference. Second, we select five classes from SD-198 [23], a benchmark dataset for skin disease classification, as another dataset for both segmentation and classification tasks. Note that due to the amount of manual work involved in annotation, we select those classes based on the number of samples in each class. The selected classes are Dermatofibroma (DF), Keratoacanthoma (KA), Pyogenic Granuloma (PG), Tinea Corporis (TC), and Tinea Faciale (TF). We choose 30 samples in each class for segmentation task, and we split them into 0.7, 0.1, and 0.2 ratio for training, validation, and testing, respectively.

Table 1 show the characteristics of these two datasets for both classification and segmentation tasks broken down by the disease type and skin tone, as calculated by the Individual Typology Angle (ITA) [24]. Specifically, we consider tan2, tan1, and dark as dark skin (ds) and others as light skin (ls). Compared to other skin tone classification schemas such as Fitzpartick scale [25], we divide ITA scores into more detailed categories (eight). One prominent observation is

Table 2: Segmentation: Performance and Fairness (margin of error reported in parenthesis)

	Method	Unet	Polar	MFSNet	ViT-Adapter	EdgeMixup
Skin	Jaccard J <sub>gap</sub>	$\begin{array}{c} 0.7053 (0.0035) \\ 0.0809 (0.0001) \end{array}$	$\begin{array}{c} 0.7126 (0.0033) \\ 0.0813 (0.0001) \end{array}$	$\begin{array}{c} 0.5877 (0.0080) \\ 0.1291 (0.0076) \end{array}$	$\begin{array}{c} 0.7027 (0.0057) \\ 0.2346 (0.0035) \end{array}$	$egin{array}{l} 0.7807(0.0031) \ 0.0379(0.0001) \end{array}$
SD-seg	Jaccard J <sub>gap</sub>	$\begin{array}{c} 0.7134 (0.0031) \\ 0.0753 (0.0001) \end{array}$	$\begin{array}{c} 0.6527 (0.0036) \\ 0.1210 (0.0003) \end{array}$	$\begin{array}{c} 0.6170 (0.0052) \\ 0.0636 (0.0033) \end{array}$	$\begin{array}{c} 0.5088 (0.0042) \\ 0.2530 (0.0021) \end{array}$	$egin{array}{l} 0.7799(0.0031) \ 0.0528(0.0001) \end{array}$

that is images are more abundant than ds images due to a disparity in the availability of ds imagery found from either public sources or from clinicians with patient consent.

## 5 Evaluation

We implement EDGEMIXUP using python 3.8 and Pytorch, and all experiments are performed using one GeForce RTX 3090 graphics card (NVIDIA).

Segmentation Evaluation. Our segmentation evaluation adopts four baselines, (i) a U-Net trained to segment skin lesions, (ii) a polar training [15] transforming images from Cartesian coordinates to polar coordinates, (iii) ViT-Adapter [16], a state-of-the-art semantic segmentation using a fine-tuned ViT model, (iv) MFSNet [17], a segmentation model with differently scaled feature maps to compute the final segmentation mask. We follow the default setting from each paper for evaluation. Our evaluation metrics include (i) Jaccard index (IoU score), which measures the similarity between a predicted mask and the manually annotated ground truth, and (ii) the gap between Jaccard values  $(J_{gap})$  to measure fairness.

Table 2 shows the performance and fairness of EDGEMIXUP and different baselines. We compare predicted masks with the manually-annotated ground truth by calculating the Jaccard index, and computing the gap for subpopulations with ls and ds (based on ITA). EDGEMIXUP, a data preprocessing method, improves the utility of lesion segmentation in terms of Jaccard index compared with all existing baselines. One reason is that EDGEMIXUP preserves skin lesion information, thus improving the segmentation quality, while attenuating markers for protected factors. Note that EDGEMIXUP iteratively improves the segmentation results. Take our Skin-seg dataset for example. We trained our baseline Unet model for three iterations, and the model utility is increased by 0.0468 on Jaccard index while the J<sub>gap</sub> between subpopulations is reduced by 0.0193.

**Classification Evaluation.** Our classification evaluation involves: (i) Adversarial Debiasing (AD) [26], (ii) DexiNed-avg, the average version of DexiNed [27] as an boundary detector used by EDGEMIXUP, and (iii) ST-Debias [22], a debiasing method augmenting data with conflicting shape and texture information. Our evaluation metrics include accuracy gap, the (Rawlsian) minimum accuracy

#### 8 Yuan et al.

Table 3: Skin disease classification and associated bias. Samples contain skin tones as a protected factor. (margin of error reported in parentheses, subpopulation reported in brackets)

	Metrics	ResNet34		Baselines	EdgeMixup (ours)			
			AD	DexiNed-avg	ST-Debias	U-Net	Mask-based	
Skin	acc acc <sub>gap</sub> acc <sub>min</sub> CAI <sub>0.5</sub> CAI <sub>0.75</sub>	88.08(3.66) 16.38(12.21) 73.33[ds] - -	$\begin{array}{c} 81.79(4.35)\\ 5.33(11.69)\\ 76.92[\mathrm{ds}]\\ 2.380\\ 6.715\end{array}$	$\begin{array}{c} 69.87(5.17)\\ 19.79(13.52)\\ 51.85[\mathrm{ds}]\\ -10.81\\ -7.110\end{array}$	76.52(5.23) 2.64(8.05) 71.12[ds] 1.090 7.415	$\begin{array}{c} 86.75(3.82)\\ 8.280(9.66)\\ 79.41[\mathrm{ds}]\\ 3.385\\ 5.743\end{array}$	86.09(3.90) <b>1.923(8.49)</b> <b>84.38[ds]</b> <b>6.233</b> <b>10.35</b>	
	AUC AUCgap AUCmin CAUCI0.5 CAUCI0.75	<b>0.977(0.02)</b> 0.039(0.07) 0.942[ds]	$\begin{array}{c} 0.956(0.02)\\ 0.009(0.05)\\ 0.955[\mathrm{ds}]\\ 0.004\\ 0.017\end{array}$	$\begin{array}{c} 0.889(0.04)\\ 0.090(0.11)\\ 0.807[\mathrm{ds}]\\ -0.069\\ -0.060\end{array}$	$\begin{array}{c} 0.933(0.03)\\ 0.035(0.04)\\ 0.910[\mathrm{ds}]\\ -0.024\\ -0.014 \end{array}$	$\begin{array}{c} 0.974(0.02)\\ 0.011(0.02)\\ 0.973[\mathrm{ds}]\\ 0.012\\ 0.020 \end{array}$	0.973(0.02) 0.01 (0.05) 0.964[ds] 0.013 0.022	
SD-sub	acc acc <sub>gap</sub> acc <sub>min</sub> CAI <sub>0.5</sub> CAI <sub>0.75</sub>	75.60(14.26) 28.12(54.98) 50.00[ds] -	$73.53(14.83) \\ 25.00(54.30) \\ 50.00[ds] \\ 0.525 \\ 1.822$	63.13(16.08) 25.21(51.20) 43.75 [ds] -4.780 -0.934	$71.73(13.01) \\18.66(13.21) \\70.59[ls] \\2.795 \\6.127$	$74.17(13.30) \\18.51(11.50) \\72.11[ls] \\4.090 \\6.850$	$76.47(11.26) \\ 15.00(9.62) \\ 75.00[ls] \\ 6.995 \\ 10.06$	
	AUC AUC <sub>gap</sub> AUC <sub>min</sub> CAUCI <sub>0.5</sub> CAUCI <sub>0.75</sub>	0.922(0.10) 0.429(0.35) 0.500[ds] -	$\begin{array}{c} 0.962(0.06)\\ 0.319(0.36)\\ 0.650[\mathrm{ds}]\\ 0.075\\ 0.092 \end{array}$	$\begin{array}{c} 0.824(0.13)\\ 0.175(0.29)\\ 0.650[\mathrm{ds}]\\ 0.078\\ 0.166\end{array}$	$\begin{array}{c} 0.941(0.08)\\ 0.255(0.31)\\ 0.711[\mathrm{ds}]\\ 0.097\\ 0.135\end{array}$	$\begin{array}{c} 0.953(0.11)\\ 0.178(0.29)\\ 0.784[\mathrm{ds}]\\ 0.140\\ 0.196\end{array}$	$\begin{array}{c} 0.970(0.06)\\ 0.170(0.29)\\ 0.800[\mathrm{ds}]\\ 0.153\\ 0.206\end{array}$	

across subpopulations, area under the receiver operating characteristic curve (AUC), and joint metrics (CAI<sub> $\alpha$ </sub> and CAUCI<sub> $\alpha$ </sub>).

Table 3 shows utility performance (acc and AUC) and fairness results (gaps of acc and AUC between ls and ds subpopulations). We here list two variants of EDGEMIXUP, and one of which, "Unet", uses the lesion edge generated by the baseline Unet model while "mask-based" implements deep-learning model involved methodology introduced in Section 3. By adding the "Unet" variant, we demonstrate here that simply applying lesion edge predicetd by the baseline Unet model, while not optimal, efficiently reduces model bias on different skin-tone samples. EDGEMIXUP outperforms SOTA approaches in balancing the model's performance and fairness, i.e., the CAI<sub> $\alpha$ </sub> and CAUCI<sub> $\alpha$ </sub> values of EDGEMIXUP are the highest compared with the vanilla ResNet34 and other baselines.

## 6 Related Work

Skin Disease Classification and Segmentation: Previous researches mainly work on improving model utility for both medical image [28] and skin lesion [29] classification. As for skin lesion segmentation tasks, few works has been proposed due to the lack of datasets with ground-truth segmentation masks. International Skin Imaging Collaboration (ISIC) hosts challenges of International Symposium on Biomedical Imaging (ISBI) [11] to encourage researches studying lesion segmentation, feature detection, and image classification. However, official datasets released, e.g., HAM10000 [10] only contains melanoma samples and all of the samples are with light skins according to our inspection using ITA scores.

**Bias Mitigation:** Researchers have addressed bias and heterogeneity in deep learning models [18, 30]. First, masking sensitive factors in imagery is shown to improve fairness in object detection and action recognition [31]. Second, adversarial debiasing operates on the principle of simultaneously training two networks with different objectives [32]. The competing two-player optimization paradigm is applied to maximizing equality of opportunity [33]. As a comparison, EDGEMIXUP is an effective preprocessing approach to debiasing when applied to skin disease particularly for Lyme-focused classification and segmentation tasks.

## 7 Conclusion

We present a simple yet novel approach to segment Lyme disease lesion, which can be further used for disease classification. The key insight is a novel data preprocessing method that utilizes edge detection and mixup to isolate and highlight skin lesions and reduce bias. EDGEMIXUP outperforms SOTAs in terms of Jaccord index for segmentation and  $CAUCI_{\alpha}$  for disease classification.

Acknowledgement. This work was supported in part by Johns Hopkins University Institute for Assured Autonomy (IAA) with grants 80052272 and 80052273, and National Science Foundation (NSF) under grants CNS18-54000. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of NSF or JHU-IAA.

## References

- Daniel SW Ting, Yong Liu, Philippe Burlina, Xinxing Xu, Neil M Bressler, and Tien Y Wong, "Ai for medical imaging goes deep," *Nature medicine*, vol. 24, no. 5, pp. 539–540, 2018.
- Philippe Burlina, Neil Joshi, William Paul, Katia D Pacheco, and Neil M Bressler, "Addressing artificial intelligence bias in retinal disease diagnostics," *Translational Vision Science and Technology*, 2020.
- Alison F. Hinckley, Neeta P. Connally, James I. Meek, Barbara J. Johnson, Melissa M. Kemperman, Katherine A. Feldman, Jennifer L. White, and Paul S. Mead, "Lyme Disease Testing by Large Commercial Laboratories in the United States," *Clinical Infectious Diseases*, 2014.
- Kiersten J Kugeler, Amy M Schwartz, Mark J Delorey, Paul S Mead, and Alison F Hinckley, "Estimating the frequency of lyme disease diagnoses, united states, 2010– 2018," *Emerging Infectious Diseases*, 2021.
- Robert B. Nadelman, "Erythema migrans," Infectious Disease Clinics of North America, 2015.

- 10 Yuan et al.
- Carrie D. Tibbles and Jonathan A. Edlow, "Does This Patient Have Erythema Migrans?," JAMA, 2007.
- Philippe M Burlina, Neil J Joshi, Elise Ng, Seth D Billings, Alison W Rebman, and John N Aucott, "Automated detection of erythema migrans and other confounding skin lesions via deep learning," *Computers in biology and medicine*, vol. 105, pp. 151–156, 2019.
- 8. Yanyang Gu, Zongyuan Ge, C Paul Bonnington, and Jun Zhou, "Progressive transfer learning and adversarial domain adaptation for cross-domain skin disease classification," *IEEE journal of biomedical and health informatics*, 2019.
- Philippe M Burlina, Neil J Joshi, Phil A Mathew, William Paul, Alison W Rebman, and John N Aucott, "Ai-based detection of erythema migrans and disambiguation against other skin lesions," *Computers in Biology and Medicine*, 2020.
- Philipp Tschandl, Cliff Rosendahl, and Harald Kittler, "The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions," *Scientific data*, 2018.
- 11. Noel Codella, Veronica Rotemberg, Philipp Tschandl, M. Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, Harald Kittler, and Allan Halpern, "Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic)," 2019.
- 12. Matthew Groh, Caleb Harris, Luis Soenksen, Felix Lau, Rachel Han, Aerin Kim, Arash Koochek, and Omar Badri, "Evaluating deep neural networks trained on clinical images in dermatology with the fitzpatrick 17k dataset," in *Proceedings of* the *IEEE/CVF CVPR*, 2021.
- Muhammad Attique Khan, Muhammad Sharif, Tallha Akram, Robertas Damaševičius, and Rytis Maskeliūnas, "Skin lesion segmentation and multiclass classification using deep learning features and improved moth flame optimization," *Diagnostics*, 2021.
- 14. Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*, 2015.
- Marin Benčević, Irena Galić, Marija Habijan, and Danilo Babin, "Training on polar image transformations improves biomedical image segmentation," *IEEE Access*, 2021.
- 16. Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao, "Vision transformer adapter for dense predictions," in *ICLR*, 2023.
- 17. Hritam Basak, Rohit Kundu, and Ram Sarkar, "Mfsnet: A multi focus segmentation network for skin lesion segmentation," *Pattern Recognition*, 2022.
- Simon Caton and Christian Haas, "Fairness in machine learning: A survey," arXiv preprint arXiv:2010.04053, 2020.
- Philippe Burlina, William Paul, Philip Mathew, Neil Joshi, Katia D Pacheco, and Neil M Bressler, "Low-shot deep learning of diabetic retinopathy with potential applications to address artificial intelligence bias in retinal diagnostics and rare ophthalmic diseases," JAMA ophthalmology, 2020.
- 20. "Edgemixup repository," https://github.com/Haolin-Yuan/EdgeMixup.
- Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell, "Mitigating unwanted biases with adversarial learning," in *Proceedings of the 2018 AAAI/ACM Conference* on AI, Ethics, and Society, 2018.
- 22. Yingwei Li, Qihang Yu, Mingxing Tan, Jieru Mei, Peng Tang, Wei Shen, Alan Yuille, and Cihang Xie, "Shape-texture debiased neural network training," in *International Conference on Learning Representations*, 2021.

- Xiaoxiao Sun, Jufeng Yang, Ming Sun, and Kai Wang, "A benchmark for automatic visual classification of clinical skin disease images," in *European Conference on Computer Vision*, 2016.
- 24. Marcus Wilkes, Caradee Y Wright, Johan L du Plessis, and Anthony Reeder, "Fitzpatrick skin type, individual typology angle, and melanin index in an african population: steps toward universally applicable skin photosensitivity assessments," JAMA dermatology, 2015.
- Thomas B. Fitzpatrick, "Soleil et peau," Journal de Médecine Esthétique (in French), 1975.
- 26. Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell, "Mitigating unwanted biases with adversarial learning," in *Proceedings of the 2018 AAAI/ACM Conference* on AI, Ethics, and Society, 2018.
- 27. X. Soria, E. Riba, and A. Sappa, "Dense extreme inception network: Towards a robust cnn model for edge detection," in *WACV*, Mar 2020.
- Zhuoning Yuan, Yan Yan, Milan Sonka, and Tianbao Yang, "Large-scale robust deep auc maximization: A new surrogate loss and empirical studies on medical image classification," in *ICCV*, 2021.
- 29. Bhuvaneshwari Shetty, Roshan Fernandes, Anisha P Rodrigues, Rajeswari Chengoden, Sweta Bhattacharya, and Kuruva Lakshmanna, "Skin lesion classification of dermoscopic images using machine learning and convolutional neural network," *Scientific Reports*, 2022.
- Haolin Yuan, Bo Hui, Yuchen Yang, Philippe Burlina, Neil Zhenqiang Gong, and Yinzhi Cao, "Addressing heterogeneity in federated learning via distributional transformation," in *ECCV*, 2022.
- Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez, "Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations," in *ICCV*, 2019.
- 32. Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein, "Adversarial training for free!," Advances in Neural Information Processing Systems, 2019.
- Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi, "Data decisions and theoretical implications when adversarially learning fair representations," arXiv preprint arXiv:1707.00075, 2017.