# Pseudo Bias-Balanced Learning for Debiased Chest X-ray Classification

Luyang Luo[1]✉, Dunyuan Xu[1], Hao Chen[2],
Tien-Tsin Wong[1], and Pheng-Ann Heng[1]

[1]Department of Computer Science and Engineering,
The Chinese University of Hong Kong, Hong Kong, China
`lyluo@cse.cuhk.edu.hk`
[2]Department of Computer Science and Engineering,
The Hong Kong University of Science and Technology, Hong Kong, China

**Abstract.** Deep learning models were frequently reported to learn from shortcuts like dataset biases. As deep learning is playing an increasingly important role in the modern healthcare system, it is of great need to combat shortcut learning in medical data as well as develop unbiased and trustworthy models. In this paper, we study the problem of developing debiased chest X-ray diagnosis models from the biased training data without knowing exactly the bias labels. We start with the observations that the imbalance of bias distribution is one of the key reasons causing shortcut learning, and the dataset biases are preferred by the model if they were easier to be learned than the intended features. Based on these observations, we proposed a novel algorithm, pseudo bias-balanced learning, which first captures and predicts per-sample bias labels via generalized cross entropy loss and then trains a debiased model using pseudo bias labels and bias-balanced softmax function. We constructed several chest X-ray datasets with various dataset bias situations and demonstrated with extensive experiments that our proposed method achieved consistent improvements over other state-of-the-art approaches.[1]

**Keywords:** Debias; Shortcut Learning; Chest X-ray

## 1   Introduction

To date, deep learning (DL) has achieved comparable or even superior performance to experts on many medical image analysis tasks [17]. Robust and trustworthy DL models are hence of greater need than ever to unleash their huge potential in solving real-world healthcare problems. However, a common trust failure of DL was frequently found where the models reach a high accuracy without learning from the intended features. For example, using backgrounds to distinguish foreground objects [19], using the gender to classify hair colors [20], or worse yet, using patients' position to determine COVID-19 pneumonia from

---

[1] Code available at https://github.com/LLYXC/PBBL.

chest X-rays [4]. Such a phenomenon is called *shortcut learning* [5], where the DL models choose unintended features, or *dataset bias*, for making decisions.

More or less, biases could be generated during the creation of the datasets [22]. As the dataset biases frequently co-occurred with the primary targets, the model might take shortcuts by learning from such spurious correlation to minimize the empirical risk over the training data. As a result, dramatic performance drops could be observed when applying the models onto other data which do not obtain the same covariate shift [13]. In the field of medical image analysis, shortcut learning has also been frequently reported, such as using hospital tokens to recognize pneumonia cases [25]; learning confounding patient and healthcare variables to identify fracture cases; relying on chest drains to classify pneumothorax case [16]; or leveraging shortcuts to determine COVID-19 patients [4]. These findings reveal that shortcut learning makes deep models less explainable and less trustworthy to doctors as well as patients, and addressing shortcut learning is a far-reaching topic for modern medical image analysis.

To combat shortcut learning and develop debiased models, a branch of previous works use data re-weighting to learn from less biased data. For instance, REPAIR [12] proposed to solve a minimax problem between the classifier parameters and dataset re-sampling weights. Group distributional robust optimization [20] prioritized worst group learning, which was also mianly implemented by data re-weighting. Yoon et al. [24] proposed to address dataset bias with a weighted loss and a dynamic data sampler. Another direction of works emphasizes learning invariance across different environments, such as invariant risk minimization [1], contrastive learning [21], and mutual information minimization [29]. However, these methods all required dataset biases to be explicitly annotated, which might be infeasible for realistic situations, especially for medical images. Recently, some approaches have made efforts to relax the dependency on explicit bias labels. Nam et al. [15] proposed to learn a debiased model by mining the high-loss samples with a highly-biased model. Lee et al. [11] further incorporated feature swapping between the biased and debiased models to augment the training samples. Yet, very few methods attempted to efficiently address shortcut learning in medical data without explicitly labeling the biases.

In this paper, we are pioneered in tackling the challenging problem of developing debiased medical image analysis models without explicit labels on the bias attributes. We first observed that the imbalance of bias distribution is one of the key causes to shortcut learning, and dataset biases would be preferred when they were easier to be learned than the intended features. We thereby proposed a novel algorithm, namely pseudo bias-balanced learning (PBBL). PBBL first develops a highly-biased model by emphasizing learning from the easier features. The biased model is then used to generate pseudo bias labels that are later utilized to train a debiased model with a bias-balanced softmax function. We constructed several chest X-ray datasets with various bias situations to evaluate the efficacy of the debiased model. We demonstrated that our method was effective and robust under all scenarios and achieved consistent improvements over other state-of-the-art approaches.

## 2   Methodology

### 2.1   Problem Statement and Study Materials

Let $X$ be the set of input data, $Y$ the set of target attributes that we want the model to learn, and $B$ the set of bias attributes that are irrelevant to the targets. Our goal is to learn a function $f : X \rightarrow Y$ that would not be affected by the dataset bias. We here built the following chest X-ray datasets for our study.

**Source-biased Pneumonia (SbP):** For the training set, we first randomly sampled 5,000 pneumonia cases from MIMIC-CXR [8] and 5,000 healthy cases (no findings) from NIH [23]. We then sampled $5,000 \times r\%$ pneumonia cases from NIH and the same amount of healthy cases from MIMIC-CXR. Here, the `data source` became the dataset bias, and `health condition` was the target to be learned. We varied $r$ to be 1, 5, and 10, which led to biased sample ratios of 99%, 95%, and 90%, respectively. We created the validation and the testing sets by equally sampling 200 and 400 images from each group (w/ or w/o pneumonia; from NIH or MIMIC-CXR), respectively. Moreover, as overcoming dataset bias could lead to better external validation performance [5], we included 400 pneumonia cases and 400 healthy cases from Padchest [2] to evaluate the generalization capability of the proposed method. Note that we converted all images to JPEG format to prevent the data format from being another dataset bias.

**Gender-biased Pneumothorax (GbP):** Previous study [10] pointed out that gender imbalance in medical datasets could lead to a biased and unfair classifier. Based on this finding, we constructed two training sets from the NIH dataset [23]: 1) **GbP-Tr1**: 800 male samples with pneumothorax, 100 male samples with no findings, 800 female samples with no findings, and 100 female samples with pneumothorax; 2)**GbP-Tr2**: 800 female samples with pneumothorax, 100 female samples with no findings, 800 male samples with no findings, and 100 male samples with pneumothorax. For validation and testing sets, we equally collected 150 and 250 samples from each group (w/ or w/o pneumothorax; male or female), respectively. Here, `gender` became a dataset bias and `health condition` was the target that the model was aimed to learn.

Following previous studies [15,11], we call a sample bias-aligned if its target and bias attributes are highly-correlated in the training set (e.g., (`pneumonia, MIMIC-CXR`) or (`healthy, NIH`) in the SbP dataset). On the contrary, a sample is said to be bias-conflicting if the target and bias attributes are dissimilar to the previous situation (e.g., (`pneumonia, NIH`) or (`healthy, MIMIC-CXR`)).

### 2.2   Bias-balanced Softmax

Our first observation is that *bias-imbalanced training data leads to a biased classifier*. Based on the SbP dataset, we trained two different settings: i) SbP with $r = 10$; ii) Bias balancing by equally sampling 500 cases from each group. The results are shown in Fig. 1a and Fig. 1b, respectively. Clearly, when the dataset is bias-imbalanced, learning bias-aligned samples were favored. On the contrary, balancing the biases mitigates shortcut learning even with less training data.
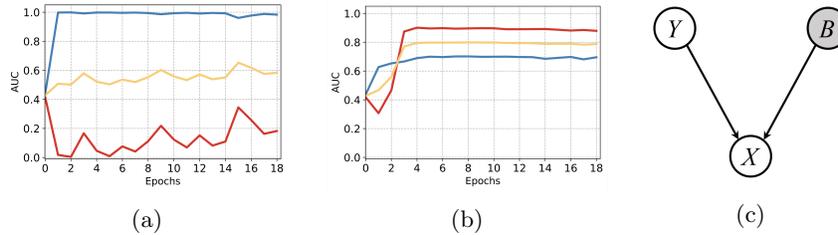
Fig. 1: We show (a) the testing results in AUC curves of a model trained on Source-biased Pneumonia dataset; (b) the testing results in AUC curves of a model trained with bias-balanced pneumonia dataset. We further show our causal assumption of data generation process in (c). Blue curves: results on bias-aligned samples; Red curves: results on bias-conflicting samples; Yellow curves: averaged results of bias-aligned AUC and bias-conflicting AUC.

For a better interpretation, we adopt the causal assumption [14] that the data $X$ is generated from both the target attributes $Y$ and the bias attributes $B$, which are independent to each other, as shown in Fig. 1c. The conditional probability $p(y = j|x)$ hence can be formalized as follows:

$$p(y = j|x, b) = \frac{p(x|y = j, b)p(y = j|b)}{p(x|b)}, \tag{1}$$

where $p(y = j|b)$ raises a distributional discrepancy between the biased training data and the ideal bias-balanced data (e.g., the testing data). Moreover, according to our experimental analysis before, the imbalance also made the model favor learning from bias-aligned samples, which finally resulted in a biased classifier. To tackle the bias-imbalance situation, let $k$ be the number of classes and $n_{j,b}$ the number of training data of target class $j$ with bias class $b$, we could derive a bias-balanced softmax [6,18] as follows:

**Theorem 1.** *(Bias-balanced softmax [6]) Assume $\phi_j = p(y = j|x, b) = \frac{p(x|y=j,b)}{p(x|b)} \cdot \frac{1}{k}$ to be the desired conditional probability of the bias-balanced dataset, and $\hat{\phi}_j = \frac{p(x|y=j,b)}{\hat{p}(x|b)} \cdot \frac{n_{j,b}}{\sum_{i=1}^{k} n_{i,b}}$ to be the conditional probability of the biased dataset. If $\phi$ can be expressed by the standard Softmax function of the logits $\eta$ generated by the model, i.e., $\phi_j = \frac{exp(\eta_j)}{\sum_{i=1}^{k} exp(\eta_i)}$, then $\hat{\phi}$ can be expressed as*

$$\hat{\phi}_j = \frac{p(y = j|b) \cdot \exp(\eta_j)}{\sum_{i=1}^{k} p(y = i|b) \cdot \exp(\eta_i)}. \tag{2}$$

Theorem 1 (proof provided in the supplementary) shows that bias-balanced softmax could well solve the distributional discrepancy between the bias-imbalanced

training set and the bias-balanced testing set. Denoting $M$ the number of training data, we obtain the bias-balanced loss for training a debiased model:

$$\mathcal{L}_{\text{BS}}(f(x), y, b) = -\frac{1}{M}\sum_{i=1}^{M} log\left(\frac{p(y=j|b)\cdot \exp(\eta_j)}{\sum_{i=1}^{k} p(y=i|b)\cdot \exp(\eta_i)}\right) \tag{3}$$

However, this loss requires estimation of the bias distribution on the training set, while comprehensively labeling all kinds of attributes would be unpractical, especially for medical data. In the next section, we elaborate on how to obtain the estimation of the bias distribution without knowing the bias labels.

### 2.3   Bias Capturing with Generalized Cross Entropy Loss

Inspired by [15], we conducted two experiments based on the Source-biased Pneumonia dataset with $r = 10$, where we set the models to classify data source (Fig. 2a) or health condition (Fig. 2b), respectively. Apparently, the model has almost no signs of fitting on the bias attribute (health condition) when it's required to distinguish data source. On the other hand, the model quickly learns the biases (data source) when set to classify pneumonia from healthy cases. From these findings, one could conclude that *dataset biases would be preferred when they were easier to be learned than the intended features.*
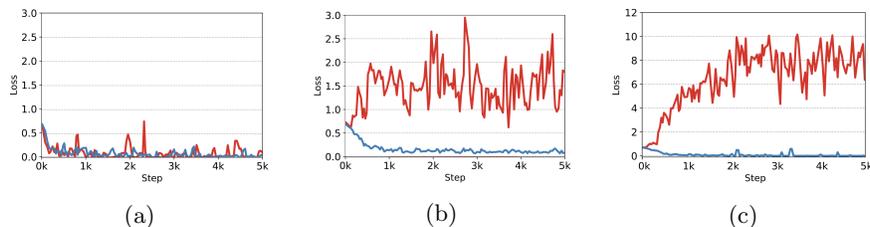


Fig. 2: Based on the SbP dataset, we show the learning curve of the vanilla model by setting the {target, bias} pair to be (a) {data source, health condition} and (b) {health condition, data source}. We also show in (c) the learning curve of a highly-biased model trained with GCE loss with the {target, bias} pair being {health condition, data source}. Blue curves: loss of bias-aligned samples; Red curves: loss of bias-conflicting samples.

Based on this observation, we could develop a model to capture the dataset bias by making it quickly fit on the easier features from the training data. Therefore, we adopt the generalized cross entropy (GCE) loss [27], which was originally proposed to address noisy labels by fitting on the easier clean data and slowly memorizing the hard noisy samples. Inheriting this idea, the GCE loss could also quickly capture easy and biased samples than the categorical cross entropy (CE) loss. Giving $f(x;\theta)$ the softmax output of the model, denoting $f_{y=j}(x;\theta)$

the probability of $x$ being classified to class $y = j$ and $\theta$ the parameters of model $f$, the GCE loss is formulated as follows:

$$\mathcal{L}_{\text{GCE}}(f(x;\theta), y = j) = \frac{1 - f_{y=j}(x;\theta)^q}{q},\tag{4}$$

where $q$ is a hyper-parameter. The gradient of GCE is $\frac{\partial \mathcal{L}_{\text{GCE}}(f(x;\theta),y=j)}{\partial \theta} = f_{y=j}(x;\theta)^q \frac{\partial \mathcal{L}_{\text{CE}}(f(x;\theta),y=j)}{\partial \theta}$ (proof provided in the supplementary), which explicitly assigns weights on the CE loss based on the agreement between model's predictions and the labels. As shown in Fig. 2c, GCE loss fits the bias-aligned samples quickly while yields much higher loss on the bias-conflicting samples.

### 2.4   Bias-balanced Learning with Pseudo Bias

With the afore discussed observations and analysis, we propose a debiasing algorithm, namely Pseudo Bias Balanced Learning. We first train a biased model $f_B(x;\theta_B)$ with the GCE loss and calculate the corresponding receiver operating characteristics (ROC) over the training set. Based on the ROC curve, we compute the sensitivity $u(\tau)$ and specificity $v(\tau)$ under each threshold $\tau$ and then assign pseudo bias labels to each sample with the following:

$$\tilde{b}(f_B(x;\theta_B)) = \begin{cases} 1, & \text{if } f_B(x;\theta_B) \geq \text{argmax}_\tau(u(\tau) + v(\tau)); \\ 0, & \text{otherwise.} \end{cases}\tag{5}$$

---

**Algorithm 1** Pseudo Bias Balanced Learning

---

**Input:** $\theta_B$, $\theta_D$, image $x$, target label $y$, numbers of iterations $T_B$, $T_D$, $N$.
**Output:** Debiased model $f_D(x;\theta_D)$.

1: Initialize $\tilde{b} = y$.
2: **for** n=1, $\cdots$, $N$ **do**
3:     Initialize network $f_B(x;\theta_B)$.
4:     **for** t=1, $\cdots$, $T_B$ **do**
5:         Update $f_B(x;\theta_B)$ with $\mathcal{L}_{\text{GCE}}(f_B(x;\theta_B),\tilde{b})$
6:     **end for**
7:     Calculate $u$, $v$, and $\tau$ over training set.
8:     Update pseudo bias labels $\tilde{b}$ with Eq. 5.
9: **end for**
10: Initialize network $f_D(x;\theta_D)$.
11: **for** t=1, $\cdots$, $T_D$ **do**
12:     Update $f_D(x;\theta_D)$ with $\mathcal{L}_{\text{BS}}(f_D(x;\theta_D),y,\tilde{b})$
13: **end for**

---

Moreover, as the biased model could also memorize the correct prediction for the hard bias-conflicting cases [26], we propose to capture and enhance the bias via iterative model training. Finally, we train our debiased model $f_D(x;\theta_D)$ based on the pseudo bias labels and the bias-balance softmax function, with different weights from $\theta_B$. The holistic approach is summarized in Algorithm 1.

## 3   Experiments

**Evaluation metrics** are the area under the ROC curve (AUC) with four criteria: i) AUC on bias-aligned samples; ii) AUC on bias-conflicting samples; iii) Average of bias-aligned AUC and bias-conflicting AUC, which we call balanced-AUC; iv) AUC on all samples. The difference between the first two metrics could reflect whether the model is biased, while the latter two metrics provide unbiased evaluations on the testing data.

   **Compared methods** included four other approaches: i) Vanilla model, which did not use any debiasing strategy and could be broadly regarded as a lower bound. ii) Group Distribution Robust Optimization (G-DRO) [20], which used the bias ground truth and could be regarded as the upper bound. G-DRO divides training data into different groups according to their targets and bias labels. It then optimized the model with priority on the worst-performing group and finally achieved robustness on every single group. As in practical scenarios, the labels for the dataset biases may not be known, we also implemented iii) Learning from Failure (LfF) [15], which developed a debiased model by weighted losses from a biased model; and iv) Disentangled Feature Augmentation (DFA) [11], which was based on LfF and further adds feature swapping and augmentation between the debiased and biased models.

   **Model training protocol** is as follows: We used the same backbone, DenseNet-121 [7] with pre-trained weights from [3], for every method. Particularly, we fixed the weights of DenseNet, replaced the final output layer with three linear layers, and used the rectified linear units as the intermediate activation function. We ran each model with three different random seeds, and reported the test results corresponding to the best validation AUC. Each model is optimized with Adam [9] for around 1,000 steps with batch size of 256 and learning rate of 1e-4. $N$ in Algorithm 1 is empirically set to 1 for SbP dataset and 2 for GbP dataset, respectively. $q$ in GCE loss is set to 0.7 as recommended in [27].

Table 1: AUC results on SbP dataset. Best results without ground truth bias labels are emphasized in **bold**. † means the method uses ground truth bias labels.

| Bias Ratio | Method | Aligned | Conflicting | Balanced | Overall | External |
|---|---|---|---|---|---|---|
| 90% | G-DRO† [20] | $70.02_{\pm2.20}$ | $89.80_{\pm0.87}$ | $79.94_{\pm0.68}$ | $80.23_{\pm0.37}$ | $90.06_{\pm0.32}$ |
| | Vanilla | $\mathbf{96.51}_{\pm0.26}$ | $31.21_{\pm3.04}$ | $63.86_{\pm1.39}$ | $69.84_{\pm1.32}$ | $71.57_{\pm0.90}$ |
| | LfF [15] | $68.57_{\pm2.16}$ | $\mathbf{87.46}_{\pm2.17}$ | $78.02_{\pm0.18}$ | $78.26_{\pm0.18}$ | $87.71_{\pm2.66}$ |
| | DFA [11] | $74.63_{\pm4.61}$ | $83.30_{\pm3.96}$ | $78.96_{\pm0.33}$ | $78.76_{\pm0.15}$ | $74.58_{\pm7.56}$ |
| | Ours | $76.82_{\pm2.80}$ | $85.75_{\pm0.32}$ | $\mathbf{80.49}_{\pm0.20}$ | $\mathbf{78.78}_{\pm3.02}$ | $\mathbf{89.96}_{\pm0.69}$ |
| 95% | G-DRO† [20] | $68.65_{\pm1.21}$ | $89.86_{\pm0.67}$ | $79.26_{\pm0.47}$ | $79.8_{\pm0.36}$ | $90.16_{\pm0.73}$ |
| | Vanilla | $\mathbf{97.91}_{\pm0.75}$ | $20.45_{\pm5.96}$ | $59.18_{\pm2.61}$ | $67.11_{\pm1.85}$ | $68.61_{\pm3.50}$ |
| | LfF [15] | $69.56_{\pm2.01}$ | $\mathbf{86.43}_{\pm1.67}$ | $77.99_{\pm0.18}$ | $\mathbf{78.28}_{\pm0.22}$ | $\mathbf{88.56}_{\pm3.37}$ |
| | DFA [11] | $69.04_{\pm4.21}$ | $84.94_{\pm2.56}$ | $76.99_{\pm0.85}$ | $77.26_{\pm0.49}$ | $76.37_{\pm3.26}$ |
| | Ours | $71.72_{\pm6.65}$ | $84.68_{\pm3.49}$ | $\mathbf{78.20}_{\pm0.20}$ | $78.04_{\pm3.46}$ | $82.65_{\pm0.40}$ |
| 99% | G-DRO† [20] | $74.30_{\pm2.28}$ | $85.18_{\pm1.26}$ | $79.74_{\pm0.55}$ | $79.71_{\pm0.40}$ | $89.87_{\pm0.64}$ |
| | Vanilla | $\mathbf{99.03}_{\pm0.95}$ | $4.93_{\pm3.68}$ | $51.98_{\pm1.60}$ | $59.21_{\pm3.76}$ | $60.79_{\pm0.98}$ |
| | LfF [15] | $77.50_{\pm11.08}$ | $64.38_{\pm8.75}$ | $70.94_{\pm1.30}$ | $71.86_{\pm1.72}$ | $73.90_{\pm4.42}$ |
| | DFA [11] | $69.33_{\pm1.74}$ | $75.48_{\pm2.61}$ | $72.40_{\pm0.48}$ | $72.49_{\pm0.45}$ | $61.67_{\pm6.86}$ |
| | Ours | $72.40_{\pm0.71}$ | $\mathbf{77.61}_{\pm0.45}$ | $\mathbf{75.00}_{\pm0.18}$ | $\mathbf{74.70}_{\pm0.14}$ | $\mathbf{78.87}_{\pm0.44}$ |

Table 2: AUC results on GbP dataset. Best results without ground truth bias labels are emphasized in **bold**. † means the method uses ground truth bias labels.

| Training | Method | Aligned | Conflicting | Balanced | Overall |
|---|---|---|---|---|---|
| **GbP-Tr1** | G-DRO[†] [20] | $85.81_{\pm0.16}$ | $83.96_{\pm0.17}$ | $84.86_{\pm0.05}$ | $84.93_{\pm0.01}$ |
| | Vanilla | $89.42_{\pm0.25}$ | $77.21_{\pm0.33}$ | $83.31_{\pm0.05}$ | $83.75_{\pm0.05}$ |
| | LfF [15] | $88.73_{\pm1.34}$ | $77.47_{\pm0.09}$ | $83.10_{\pm0.64}$ | $83.46_{\pm0.71}$ |
| | DFA [11] | $86.12_{\pm0.46}$ | $\mathbf{77.92_{\pm0.23}}$ | $82.02_{\pm0.31}$ | $82.23_{\pm0.30}$ |
| | Ours | $\mathbf{90.17_{\pm0.42}}$ | $77.07_{\pm1.73}$ | $\mathbf{83.62_{\pm0.68}}$ | $\mathbf{84.13_{\pm0.56}}$ |
| **GbP-Tr2** | G-DRO[†] [20] | $83.76_{\pm1.59}$ | $85.14_{\pm0.31}$ | $84.45_{\pm0.65}$ | $84.42_{\pm0.61}$ |
| | Vanilla | $\mathbf{89.39_{\pm0.85}}$ | $76.13_{\pm0.93}$ | $82.76_{\pm0.78}$ | $82.93_{\pm0.78}$ |
| | LfF [15] | $87.25_{\pm0.62}$ | $79.07_{\pm0.96}$ | $83.16_{\pm0.45}$ | $83.19_{\pm0.44}$ |
| | DFA [11] | $80.44_{\pm0.58}$ | $\mathbf{85.51_{\pm0.57}}$ | $82.98_{\pm0.19}$ | $83.09_{\pm0.21}$ |
| | Ours | $86.34_{\pm0.64}$ | $81.69_{\pm2.67}$ | $\mathbf{84.02_{\pm1.01}}$ | $\mathbf{84.03_{\pm0.97}}$ |

**Quantitative results on Source-biased Pneumonia dataset** are reported in Table 1. With the increasing of bias ratio, the vanilla model became more and more biased and severe decreases in balanced-AUC and overall-AUC was observed. All other methods also showed decreases on the two metrics, while G-DRO shows quite robust performance under all situations. Meanwhile, our method achieved consistent improvement over the compared approaches under most of the situations, demonstrating its effectiveness in debiasing. Interestingly, the change of external testing performance appeared to be in line with the change of the balanced-AUC and overall AUC, which further revealed that overcoming shortcut learning improves the model's generalization capability. These findings demonstrated our method's effectiveness in solving shortcut learning, with potential in robustness and trustworthiness for real-world clinic usage.

**Quantitative results on Gender-biased Pneumothorax dataset** are reported in Table 2. By the performance of the vanilla model, gender bias may not affect the performance as severely as data source bias, but it could lead to serious fairness issues. We observed that G-DRO showed robust performance on the two different training sets. Among approaches that do not use ground truth
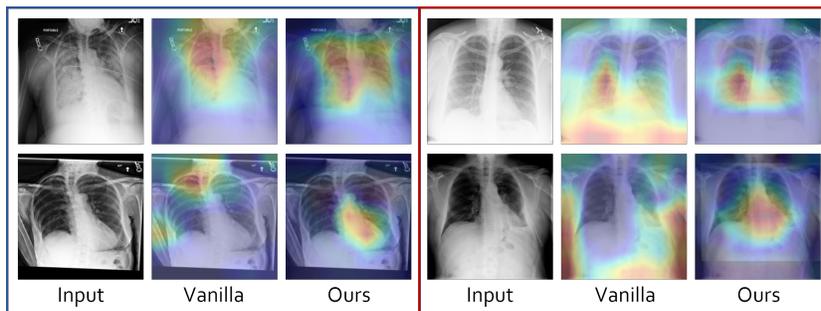


Fig. 3: Class activation map [28] generated from vanilla model and our method. Samples are from the SbP dataset (in blue box) and the GbP dataset (in red box), respectively.

bias labels, our proposed method achieved consistent improvement over others with the two different training sets. The results also showed the potential of our method in developing fair and trustworthy diagnosis models.

**Qualitative results** were visualized by class activation map [28], as shown in Fig. 3. It can be observed that vanilla model might look for evidence outside the lung regions, while our method could more correctly focus on the lung regions.

## 4    Conclusion

In this paper, we studied the causes and solutions for shortcut learning in medical image analysis, with chest X-ray as an example. We showed that shortcut learning occurs when the bias distribution is imbalanced, and the dataset bias is preferred when it is easier to be learned than the intended features. Based on these findings, we proposed a novel pseudo bias balanced learning algorithm to develop a debiased model without explicit labeling on the bias attribute. We also constructed several challenging debiasing datasets from public-available data. Extensive experiments demonstrated that our method overcame shortcut learning and achieved consistent improvements over other state-of-the-art methods under different scenarios, showing promising potential in developing robust, fair, and trustworthy diagnosis models.

## References

1. Arjovsky, M., Bottou, L., Gulrajani, I., Lopez-Paz, D.: Invariant risk minimization. arXiv preprint arXiv:1907.02893 (2019) 2

2. Bustos, A., Pertusa, A., Salinas, J.M., de la Iglesia-Vayá, M.: Padchest: A large chest x-ray image dataset with multi-label annotated reports. Medical image analysis **66**, 101797 (2020) 3

3. Cohen, J.P., Viviano, J.D., Bertin, P., Morrison, P., Torabian, P., Guarrera, M., Lungren, M.P., Chaudhari, A., Brooks, R., Hashir, M., et al.: Torchxrayvision: A library of chest x-ray datasets and models. arXiv preprint arXiv:2111.00595 (2021) 7

4. DeGrave, A.J., Janizek, J.D., Lee, S.I.: Ai for radiographic covid-19 detection selects shortcuts over signal. Nature Machine Intelligence **3**(7), 610–619 (2021) 2

5. Geirhos, R., Jacobsen, J.H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., Wichmann, F.A.: Shortcut learning in deep neural networks. Nature Machine Intelligence **2**(11), 665–673 (2020) 2, 3

6. Hong, Y., Yang, E.: Unbiased classification through bias-contrastive and bias-balanced learning. Advances in Neural Information Processing Systems **34** (2021) 4, 12

7. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 4700–4708 (2017) 7

8.  Johnson, A.E., Pollard, T.J., Berkowitz, S.J., Greenbaum, N.R., Lungren, M.P., Deng, C.y., Mark, R.G., Horng, S.: Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. Scientific data **6**(1), 1–8 (2019) 3

9.  Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Proc. Int. Conf. Learn. Representations (2015) 7

10. Larrazabal, A.J., Nieto, N., Peterson, V., Milone, D.H., Ferrante, E.: Gender imbalance in medical imaging datasets produces biased classifiers for computer-aided diagnosis. Proceedings of the National Academy of Sciences **117**(23), 12592–12594 (2020) 3

11. Lee, J., Kim, E., Lee, J., Lee, J., Choo, J.: Learning debiased representation via disentangled feature augmentation. Advances in Neural Information Processing Systems **34** (2021) 2, 3, 7, 8

12. Li, Y., Vasconcelos, N.: Repair: Removing representation bias by dataset resampling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9572–9581 (2019) 2

13. Luo, L., Chen, H., Xiao, Y., Zhou, Y., Wang, X., Vardhanabhuti, V., Wu, M., Heng, P.A.: Rethinking annotation granularity for overcoming deep shortcut learning: A retrospective study on chest radiographs. arXiv preprint arXiv:2104.10553 (2021) 2

14. Mitrovic, J., McWilliams, B., Walker, J.C., Buesing, L.H., Blundell, C.: Representation learning via invariant causal mechanisms. In: International Conference on Learning Representations (2020) 4

15. Nam, J., Cha, H., Ahn, S.S., Lee, J., Shin, J.: Learning from failure: De-biasing classifier from biased classifier. Advances in Neural Information Processing Systems **33** (2020) 2, 3, 5, 7, 8

16. Oakden-Rayner, L., Dunnmon, J., Carneiro, G., Ré, C.: Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In: Proceedings of the ACM conference on health, inference, and learning. pp. 151–159 (2020) 2

17. Rajpurkar, P., Chen, E., Banerjee, O., Topol, E.J.: Ai in health and medicine. Nature Medicine pp. 1–8 (2022) 1

18. Ren, J., Yu, C., Ma, X., Zhao, H., Yi, S., et al.: Balanced meta-softmax for long-tailed visual recognition. Advances in neural information processing systems **33**, 4175–4186 (2020) 4, 12

19. Ribeiro, M.T., Singh, S., Guestrin, C.: " why should i trust you?" explaining the predictions of any classifier. In: Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. pp. 1135–1144 (2016) 1

20. Sagawa, S., Koh, P.W., Hashimoto, T.B., Liang, P.: Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In: International Conference on Learning Representations (2020) 1, 2, 7, 8

21. Tartaglione, E., Barbano, C.A., Grangetto, M.: End: Entangling and disentangling deep representations for bias correction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13508–13517 (2021) 2

22. Torralba, A., Efros, A.A.: Unbiased look at dataset bias. In: CVPR 2011. pp. 1521–1528. IEEE (2011) 2

23. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classi-

fication and localization of common thorax diseases. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2097–2106 (2017) 3

24. Yoon, C., Hamarneh, G., Garbi, R.: Generalizable feature learning in the presence of data bias and domain class imbalance with application to skin lesion classification. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 365–373. Springer (2019) 2

25. Zech, J.R., Badgeley, M.A., Liu, M., Costa, A.B., Titano, J.J., Oermann, E.K.: Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. PLoS medicine **15**(11), e1002683 (2018) 2

26. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning requires rethinking generalization. In: ICLR (2017) 6

27. Zhang, Z., Sabuncu, M.: Generalized cross entropy loss for training deep neural networks with noisy labels. Advances in neural information processing systems **31** (2018) 5, 7, 13

28. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2921–2929 (2016) 8, 9

29. Zhu, W., Zheng, H., Liao, H., Li, W., Luo, J.: Learning bias-invariant representation by cross-sample mutual information minimization. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15002–15012 (2021) 2

## 5   Supplementary

**Proof of Theorem 1** is provided following [6,18] for better reference. The exponential family parameterization of the multinomial distribution provides the standard Softmax function as the *canonical response function* as follows:

$$\phi_j = \frac{\exp(\eta_j)}{\sum_{i=1}^{k} \exp(\eta_i)} \tag{6}$$

and the *canonical link function* as:

$$\eta_j = \log(\frac{\phi_j}{\phi_k}) \tag{7}$$

By adding $-\log(\phi_j/\hat{\phi}_j)$ to both sides of Eq. 7, we have:

$$\eta_j - \log(\frac{\phi_j}{\hat{\phi}_j}) = \log(\frac{\phi_j}{\phi_k}) - \log(\frac{\phi_j}{\hat{\phi}_j}) = \log(\frac{\hat{\phi}_j}{\phi_k}), \tag{8}$$

from which we further have:

$$\phi_k \exp(\eta_j - \log(\frac{\phi_j}{\hat{\phi}_j})) = \hat{\phi}_j \tag{9}$$

$$\phi_k \sum_{i=1}^{k} \exp(\eta_i - \log(\frac{\phi_i}{\hat{\phi}_i})) = \sum_{i=1}^{k} \hat{\phi}_i = 1 \tag{10}$$

$$\phi_k = 1/\sum_{i=1}^{k} \exp(\eta_i - \log(\frac{\phi_i}{\hat{\phi}_i})) \tag{11}$$

Substitute Eq. 11 back to Eq. 9, we could have:

$$\hat{\phi}_j = \phi_k \exp(\eta_j - \log(\frac{\phi_j}{\hat{\phi}_j})) = \frac{\exp(\eta_j - \log(\frac{\phi_j}{\hat{\phi}_j}))}{\sum_{i=1}^{k} \exp(\eta_i - \log(\frac{\phi_i}{\hat{\phi}_i}))} \tag{12}$$

We recall that

$$\phi_j = p(y=j|x,b) = \frac{p(x|y=j,b)}{p(x|b)} \cdot \frac{1}{k}; \; \hat{\phi}_j = \frac{p(x|y=j,b)}{\hat{p}(x|b)} \cdot \frac{n_{j,b}}{\sum_{i=1}^{k} n_{i,b}} \tag{13}$$

Hence,

$$\log(\frac{\phi_j}{\hat{\phi}_j}) = \log(\frac{\sum_{i=1}^{k} n_{i,b}}{k n_{j,b}}) + \log(\frac{\hat{p}(x|b)}{p(x|b)}) \tag{14}$$

For simplicity, we let $n_b = \sum_{i=1}^{k} n_{j,b}$ to be the number of samples obtaining the bias label as $b$. Finally, by substituting Eq. 14 back to Eq. 12, we have

$$\begin{aligned} \hat{\phi}_j &= \frac{\exp(\eta_j - \log\frac{n_b}{k n_{j,b}} - \log\frac{\hat{p}(x|b)}{p(x|b)})}{\sum_{i=1}^{k} \exp(\eta_i - \log\frac{n_b}{k n_{i,b}} - \log\frac{\hat{p}(x|b)}{p(x|b)})} \\ &= \frac{\frac{n_{j,b}}{n_b} \cdot \exp(\eta_j)}{\sum_{i=1}^{k} \frac{n_{i,b}}{n_b} \cdot \exp(\eta_i)} = \frac{p(y=j|b) \cdot \exp(\eta_j)}{\sum_{i=1}^{k} p(y=i|b) \cdot \exp(\eta_i)}. \end{aligned} \tag{15}$$

**Gradient of Generalized Cross Entropy Loss [27]:** The form of the GCE loss is as follows:

$$\mathcal{L}_{\text{GCE}}(f(x;\theta), y = j) = \frac{1 - f_{y=j}(x;\theta)^q}{q},$$  (16)

Hence, the gradient is:

$$\frac{\partial \mathcal{L}_{\text{GCE}}(f(x;\theta), y = j)}{\partial \theta} = -f_{y=j}(x;\theta)^{q-1} \frac{\partial f_{y=j}(x;\theta)}{\partial \theta}$$  (17)

Recall that the form of conventional cross entropy loss is $\mathcal{L}_{\text{CE}}(f(x;\theta), y = j) = -\log(f_{y=j}(x;\theta))$, hence

$$\frac{\partial \mathcal{L}_{\text{CE}}(f(x;\theta), y = j)}{\partial \theta} = -f_{y=j}(x;\theta)^{-1} \frac{\partial f_{y=j}(x;\theta)}{\partial \theta}$$  (18)

Therefore,

$$\frac{\partial \mathcal{L}_{\text{GCE}}(f(x;\theta), y = j)}{\partial \theta} = f_{y=j}(x;\theta)^q \frac{\partial \mathcal{L}_{\text{CE}}(f(x;\theta), y = j)}{\partial \theta}$$  (19)