

DEMONSTRATING THE RISK OF IMBALANCED DATASETS IN CHEST X-RAY IMAGE-BASED DIAGNOSTICS BY PROTOTYPICAL RELEVANCE PROPAGATION

Srishti Gautam*, Marina M.-C. Höhne[†]*, Stine Hansen*, Robert Jenssen* and Michael Kampffmeyer*

* UiT The Arctic University of Norway, Tromsø, Norway

[†]Technical University of Berlin, Berlin, Germany

ABSTRACT

The recent trend of integrating multi-source Chest X-Ray datasets to improve automated diagnostics raises concerns that models learn to exploit source-specific correlations to improve performance by recognizing the source domain of an image rather than the medical pathology. We hypothesize that this effect is enforced by and leverages label-imbalance across the source domains, i.e., prevalence of a disease corresponding to a source. Therefore, in this work, we perform a thorough study of the effect of label-imbalance in multi-source training for the task of pneumonia detection on the widely used ChestX-ray14 and CheXpert datasets. The results highlight and stress the importance of using more faithful and transparent self-explaining models for automated diagnosis, thus enabling the inherent detection of spurious learning. They further illustrate that this undesirable effect of learning spurious correlations can be reduced considerably when ensuring label-balanced source domain datasets.

Index Terms— Chest X-Ray, Self-Explaining Models, Explainable AI, Spurious Learning, Artifact detection.

1. INTRODUCTION

Current approaches for computer-aided diagnosis using Chest X-Ray images and deep learning tackle the lack of labeled data by leveraging data from multiple sources to achieve state-of-the-art performance [1]. However, the validity of this approach has recently been questioned by illustrating that models trained on datasets where each source exclusively contains labeled samples from a single class can directly learn the source peculiarities to solve the task [2].

In this work, we illustrate that this behavior goes far beyond this extreme setting of source-based label exclusiveness where the models are prone to relying on spurious correlations even in the presence of minor imbalances of disease prevalence across the sources. Specifically, the models tend to pick up on the textual image annotations present in the Chest X-Ray images, which may include metadata such as orien-

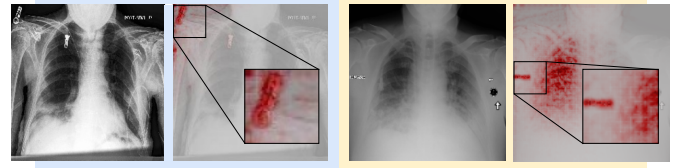


Fig. 1. Heatmaps of models for 90% (blue) and 60% (yellow) label-imbalance demonstrating spurious learning. With more imbalance, the reliance on the source annotations increases.

tation or timestamp as well as information about patients, wards, and hospitals [3]. This can lead to falsification of performance statistics as the model appears to be working well when in reality failing to capture class-related pathology-based characteristics. Our hypothesis is validated by performing a thorough analysis on the combination of the two commonly used Chest X-Ray datasets, ChestX-Ray14 [4] and CheXpert [5] for the scenario of pneumonia detection, thus simulating two different sources of X-Ray images. We deliberately introduce a gradual imbalance in the prevalence of pneumonia images from one hospital system to assess the behavior of the model. Experimental results demonstrate that the model learns source related text-annotations (see Fig. 1).

This unanticipated behavior can go unnoticed with black-box models, thus advocating the use of explainable AI. We, therefore, illustrate how these spurious correlations can be detected with the help of explainable methods. In particular, we rely on a self-explaining approach that can inherently explain the model’s underlying decision strategies without relying on post-hoc approaches, thus generating more faithful explanations [6]. We leverage Prototypical Relevance Propagation (PRP) [7], a model-aware extension of the self-explaining model ProtoPNet [8] that provides more spatially accurate and high-resolution prototypical explanations, to further support our hypothesis.

With this work, we contribute to sharpen the awareness for the use of label-balanced multi-source datasets as well as the importance of the use of self-explanatory models for

Table 1. Controlled test configurations for assessing behavior of recognizing source annotations. P and NP refer to Pneumonia and Non-Pneumonia class, respectively. The percentages denote the amount of images selected from H1 and H2.

Name	Setup		Hypothesis
Test-100H1		H1 H2	If the model relies only on the source annotations, the accuracy of this is expected to be 100% for 100H1-0H2 label-imbalance and 0% for 0H1-100H2.
	P	100% 0%	
	NP	0% 100%	
Test-100H2		H1 H2	The behavior of this setting is expected to be the opposite to that of Test-100H1.
	P	0% 100%	
	NP	100% 0%	
Test-50-50		H1 H2	The accuracy of this test set for different imbalances will indicate the learning of real disease-specific features.
	P	50% 50%	
	NP	50% 50%	

computer-aided diagnostics.

2. LABEL-IMBALANCE ANALYSIS SETUP

For illustrating the validity of our hypothesis, we first describe our setup for generating label-imbalanced multi-source datasets followed by a description of the self-explainable model that is leveraged in this study.

2.1. Datasets

To consider a controlled setting, we take a balanced subset of classes with equal number of images from both the ChestX-ray14 and CheXpert datasets for Pneumonia detection.

The NIH **ChestX-ray14** dataset consists of 112,120 frontal-view X-Ray images from 14 classes [4]. We split the data patient-wise into 80% training, 10% validation and 10% test. For our controlled setting, we first select the Pneumonia class consisting of 1099 training images. We then sample 1099 images for the negative class randomly from the remaining 13 classes to ensure a balanced dataset and to remove the additional effect of class imbalance. Similarly, this results in 374 validation and 290 total test images. We denote the images from this hospital system as H1.

CheXpert is a large public dataset consisting of 224,316 chest radiographs of 65,240 patients consisting of 14 labels. [5]. Training and validation splits for this dataset are available. We separate the training split patient-wise into train and validation data and use the dataset’s validation data for testing. To ensure two balanced datasets, we again sample 1099 images from the Pneumonia class and 1099 images from the negative class. We refer to the images from this hospital system as H2.

2.1.1. Source induced label-imbalance

In order to analyze the effect of label-imbalance across two different hospital systems (H1 and H2), we combine the two

datasets such that they vary in their composition of Pneumonia and Non-Pneumonia images obtained from H1 and H2. The substitution of the datasets is denoted as $xH1$ and $yH2$, where $x \in [0, 100]$ denotes the percentage of Pneumonia training images taken from hospital system H1 and $y = 100 - x$ denotes the percentage of Pneumonia images taken from hospital system H2. In total, we create 11 datasets and accordingly train 11 models, where x takes a value in the range of 0-100 with an interval of 10. The Non-Pneumonia images are selected such as to always maintain the class balance, i.e. $(100 - x)\%$ from H1 and $(100 - y)\%$ from H2. For example, in the case of 0H1-100H2, the training data consists of 0% Pneumonia images from H1, 100% (1099) Pneumonia images from H2, 100% (1099) Non-Pneumonia images from H1 and 0% Non-Pneumonia images from H2. The validation data is selected following the same strategy of $xH1-yH2$. For testing, three configurations are used for all models to assess their behavior of recognizing source annotations (see Table 1).

2.2. Self-explainable method: PRP

Considering the ability of black box models to learn spurious correlations [9, 10], explainable AI is essential for automated medical image diagnosis. Additionally, instead of explaining the black box models post-hoc, transparent self-explainable methods, which are capable of generating real-time explanations of the underlying decision process, can prove to be more faithful [6]. In this work, we leverage the recent method of PRP [7] to obtain class-based prototypical explanation maps.

PRP builds on ProtoPNet [8], a self-explaining model, and consists of a class-specific prototype layer inserted between the convolutional output and the final fully connected layer. The convolutional output is denoted as $\mathbf{z} \in \mathcal{R}^{H \times W \times D}$, where H , W and D are the height, width and depth of \mathbf{z} . The prototype layer consists of a fixed number of prototypes per class, $\mathbf{P} = \{\mathbf{p}_m\}_{m=1}^N$ where N are the total number of prototypes, each having a shape of $1 \times 1 \times D$. These are replaced by the closest training image patch features during training, thus representing each class by actual training image patches. L_2 similarities between the prototypes \mathbf{P} and patches of the convolutional output, $\tilde{\mathbf{z}} \in \text{patches}(\mathbf{z})$, of the input image are then computed to generate prototypical activation maps, $\mathbf{A} = \{\mathbf{a}_m\}_{m=1}^N$. This is followed by max pooling on the activation maps to generate the corresponding similarity scores $\mathbf{S} = \{s_m\}_{m=1}^N$. The network is trained in 3 steps: 1) training the whole network end to end, 2) projecting the prototypes to maintain explainability, i.e., replacing the prototypes by the convolutional output patch from the nearest training image of same class, 3) training the last layer.

PRP, unlike ProtoPNet [8], does not perform visualization of relevant areas in the input via a model-agnostic bi-linear upsampling of the activation maps to the input size, but proposes a model-aware approach inspired by LRP [11]. This leads to more faithful, higher resolution and spatially precise

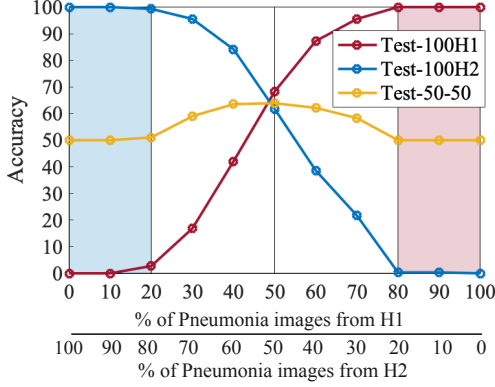


Fig. 2. Accuracy for Test-100H1, Test-100H2 and Test-50-50 for different imbalances of Pneumonia images in training datasets based on the hospital systems. 11 models are trained by combining xH1 and yH2 images, where x (percentage of Pneumonia images from H1) is in the range of 0-100 (as shown on x-axis), with an interval of 10 and $y=1-x$ is the percentage of Pneumonia images from H2.

explanation maps by taking into account the network’s structure and weights. Following PRP [7], the relevance is distributed layer by layer to the input pixels, starting from similarity scores (S).

3. EXPERIMENTS AND RESULTS

We train 11 binary classification models with the architecture in [7] for different dataset compositions as described in Section 2.1.1. Following [8], the number of prototypes for each class is fixed to 10 and ResNet34 is used as the backbone. The network is trained for 50 epochs with a projection of prototypes after every 10 epochs, followed by training the last layer for 20 epochs. The model selection is performed based on the best validation accuracy.

The test accuracies for Test-100H1, Test-100H2 and Test-50-50 for all models from 0H1-100H2 to 100H1-0H2 are shown in Fig. 2. The highest accuracy of 63.89% for the Test-50-50 (yellow) is achieved by the model trained on 50H1-50H2. Training even on a slight label-imbalanced dataset, we can observe a significant decrease in the accuracy for Test-50-50. As we move from the center of the graph i.e., 50H1-50H2 towards the left, the percentage of Pneumonia images from H2 increases and H1 decreases. Consequently, the accuracy for Test-100H1 decreases and Test-100H2 increases, strengthening our hypothesis of cheating by the model by exploiting source information. In the blue shaded region on the left, the accuracy for Test-100H2 is almost 100%, Test-100H1 is near 0% and Test-50-50 is 50%. This indicates that in the case of extreme label-imbalance, the model is only using the source related annotations for achieving better performance, thus acting as a hospital instead of a pneumonia detector. The opposite observations can be made

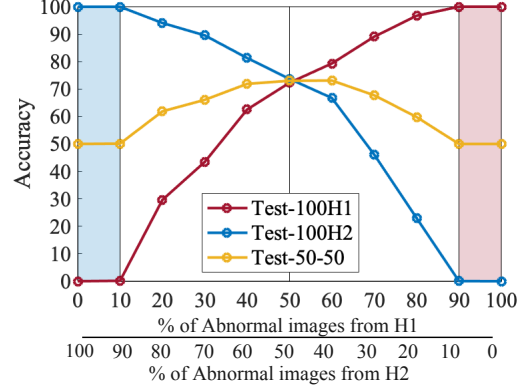


Fig. 3. Accuracy for Test-100H1, Test-100H2 and Test-50-50 for gradual imbalance of Abnormal images based on the hospital systems. Learning of source-specific annotations can be observed while inducing even a slight label-imbalance.

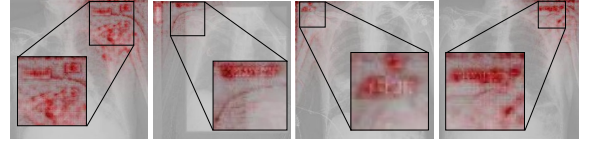


Fig. 4. Visualization of the PRP maps for the top four activated images (having the highest similarity score) by the Non-Pneumonia class prototype for 90H1-10H2 shown in Fig. 1(blue). Note, this prototype is able to capture several different kinds of annotations from hospital system H2.

when moving from the middle to the right of the graph, as the percentage of Pneumonia images are increasing from H1. The red region on the right mirrors the blue region, indicating the learning of only hospital based annotations by the models.

Considering that Pneumonia detection is a difficult problem, we also repeat the same experiments for Abnormality detection. For this, we select data from the “No Finding” category (absence of all pathologies) and data from the remaining 13 categories and consider them as the “normal” and “abnormal” class, respectively. We again gradually induce label-imbalance by varying the percentage of abnormal data coming from H1 and H2 and train 11 models for this scenario. We follow the same setting for the test configurations (see Table 1) where the percentages now correspond to normal and abnormal class data. The results are shown in Fig. 3, where the accuracy of the Test-50-50 goes up to 73.01% when using the label-balanced training data 50H1-50H2. As the imbalance increases i.e., moving to the left or right of the graph from the middle, the models start behaving as a hospital detector. Although, the effect is less severe than for the more difficult problem of Pneumonia detection, hospital detection is still observed in all cases except for the label balanced dataset of 50H1-50H2. This further stresses the importance of using label balance data for multi-source data analysis.

To demonstrate the significance of self-explaining models

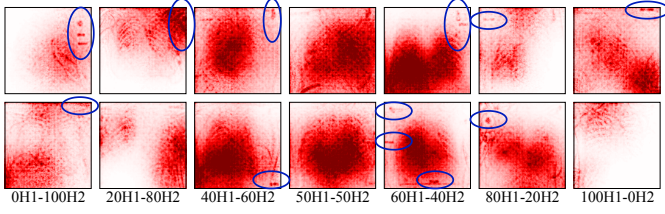


Fig. 5. Global PRP maps for Non-Pneumonia class in row 1 and Pneumonia class in row 2 for different models with learnt annotations marked in blue. As the imbalance in the hospital decreases, the models focus more on the center of the image and less on the annotations.

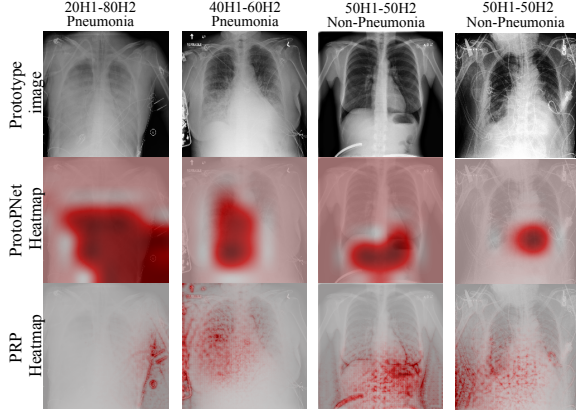


Fig. 6. Other artifacts captured by the prototypes. The rows represent the original image where the prototype comes from, ProtoPNet explanation heatmap and PRP map, respectively. Model names and prototype class are mentioned on the top.

for detecting spurious artifact learning, for a prototype from the Non-Pneumonia class of 90H1-10H2 model (Fig. 1, blue), we visualize the top 4 activated training images by this prototype in Fig. 4. The PRP maps, overlayed on the images, clearly indicate the activation of various kinds of source annotations, which are captured by the prototype, as shown in the zoomed regions.

To visualize the aggregate information learned for a class, for each model, the PRP maps for all unique prototypes learned for both, Pneumonia and Non-Pneumonia class are superimposed and shown in Fig. 5, thus representing the class-specific global prototypes. The annotations learned by the models are circled in blue. Interestingly, as we move from 50H1-50H2 to even a slight imbalance of 60-40, the models start capturing the textual annotations. From left to right in Fig. 5, it can be observed that the models start focusing more in the center of the image, trying to capture the real disease-specific features when the imbalance decreases, especially in the case of 50H1-50H2, where the global PRP maps strongly highlight the center of the images for both the classes.

In Fig. 6, we visualize the PRP maps for various prototypes capturing, among others, spurious information in the images. From the PRP maps we can observe that the model is

Table 2. AI and AD for similarity scores of predicted class prototypes on corresponding test sets for the different models i.e., Test-100H2, Test-50 and Test-100H2 for 0H1-100H2, 50H1-50H2 and 100H1-0H2, respectively. Lower A.D. and Higher A.I suggests better performance.

	A.D.		A.I.	
	ProtoPNet	PRP	ProtoPNet	PRP
0H1-100H2	13.07	9.70	38.44	43.58
50H1-50H2	52.58	48.37	20.06	22.08
100H1-0H2	12.47	12.92	45.65	48.10

capturing artifactual medical instruments, such as chest tubes, drips, and glucose bottles. It is interesting to note that in the label-balanced case, even if the reliance on source-specific artifacts is reduced, the model can still capture disease-specific artifacts in addition to the pathology features to achieve better performance (Fig. 6 column 3 and 4). An example of this can be the prevalence of chest tubes in Pneumothorax class [10], which are being captured by Non-Pneumonia class prototype in column 4 of Fig. 6, thus inaccurately diagnosing the absence of chest tubes as Pneumonia. These observations further stress on the importance of self-explainable models even for the label-balanced datasets. Additionally, to visualize the strengths of PRP over ProtoPNet, we also show the corresponding ProtoPNet heatmaps in Fig. 6 in the middle row. As can be seen the PRP maps are more precise as opposed to the inconclusive and coarse ProtoPNet explanations.

For quantifying the faithfulness of PRP maps over ProtoPNet heatmaps, we calculate the Average Drop (A.D.) and Average Increase (A.I.) with respect to similarity scores corresponding to the predicted class prototypes. Following [12], we mask out the 50% least activated pixels in the heatmaps replacing them with random uniformly sampled values. A.D is then expressed as $\sum_{m=1}^n \sum_{i=1}^K \frac{\max(0, s_m(i) - o_m(i))}{s_m(i)} \times \frac{100}{n \times K}$, where $s_m(i)$ is the similarity score for prototype m and image i , $o_m(i)$ is the output similarity score for the masked image and n, K are the number of prototypes for the predicted class and total number of images, respectively. A.I is expressed as $\sum_{m=1}^n \sum_{i=1}^K \mathbb{1}[s_m(i) < o_m(i)] \times \frac{100}{n \times K}$, where $\mathbb{1}[\cdot]$ is the Iverson bracket indicator function that returns 1 when the condition is true. Table 2 shows the values for A.I. and A.D. for both ProtoPNet and PRP maps, averaged over all images in the corresponding test sets and prototypes for models 0H1-100H2, 50H1-50H2 and 100H1-0H2 with test sets Test-100H2, Test-50-50 and Test-100H1, respectively. The results demonstrate that PRP performs consistently better or comparable to ProtoPNet for generating more faithful explanations.

4. CONCLUSION

Multi-source applicability of black box deep learning models remains questionable. In this work, we demonstrate that the models are prone to learning spurious correlations in terms

of textual annotations for Chest X-Ray image analysis in the presence of source induced label-imbalances. Even with a slight imbalance, the models are inclined to cheat and act as a hospital detector instead of the disease detector. Consequently, we recommend to ensure label-balancing while using multi-source datasets for efficient clinical deployment. Further, using a self-explainable method of PRP, we highlight the importance of using more transparent self-explainable models for real-time detection of spurious learning.

Compliance with Ethical Standards: No ethical approval was required because of the retrospective use of open source datasets.

5. REFERENCES

- [1] Hong-Yu Zhou, Shuang Yu, Cheng Bian, Yifan Hu, Kai Ma, and Yefeng Zheng, “Comparing to learn: Surpassing imagenet pretraining on radiographs by comparing image representations,” in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*.
- [2] Alex J DeGrave, Joseph D Janizek, and Su-In Lee, “AI for radiographic COVID-19 detection selects shortcuts over signal,” *Nature Machine Intelligence*, vol. 3, no. 7, pp. 610–619, 2021.
- [3] Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-Ying Deng, Roger G. Mark, and Steven Horng, “Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports,” *Scientific Data*, vol. 6, no. 1, 2019.
- [4] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M. Summers, “Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [5] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng, “Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison,” *arXiv*, 2019.
- [6] Cynthia Rudin, “Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead,” *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.
- [7] Srishti Gautam, Marina M. C. Höhne, Stine Hansen, Robert Jenssen, and Michael Kampffmeyer, “This looks more like that: Enhancing self-explaining models by prototypical relevance propagation,” *arXiv*, 2021.
- [8] Alina Barnett Jonathan Su Cynthia Rudin Chaofan Chen, Oscar Li, “This looks like that: Deep learning for interpretable image recognition,” in *Proceedings of Neural Information Processing Systems (NeurIPS)*, 2019.
- [9] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller, “Unmasking clever hans predictors and assessing what machines really learn,” *Nature Communications*, vol. 10, no. 1, Dec.
- [10] John R. Zech, Marcus A. Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and Eric Karl Oermann, “Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study,” *PLOS Medicine*, vol. 15, no. 11, pp. 1–17, 11 2018.
- [11] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek, “On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation,” *PLOS ONE*, vol. 10, no. 7, pp. 1–46, 07 2015.
- [12] Jeong Ryong Lee, Sewon Kim, Inyong Park, Taejoon Eo, and Dosik Hwang, “Relevance-cam: Your model already knows where to look,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 14944–14953.