# Slicing Through Bias: Explaining Performance Gaps in Medical Image Analysis using Slice Discovery Methods

Vincent Olesen[1], Nina Weng[1], Aasa Feragen[1], and Eike Petersen[1,2]

[1] Technical University of Denmark, Kongens Lyngby, Denmark
{ninwe, afhar, ewipe}@dtu.dk
[2] Fraunhofer Institute for Digital Medicine MEVIS, Bremen, Germany

**Abstract.** Machine learning models have achieved high overall accuracy in medical image analysis. However, performance disparities on specific patient groups pose challenges to their clinical utility, safety, and fairness. This can affect known patient groups – such as those based on sex, age, or disease subtype – as well as previously unknown and unlabeled groups. Furthermore, the root cause of such observed performance disparities is often challenging to uncover, hindering mitigation efforts. In this paper, to address these issues, we leverage Slice Discovery Methods (SDMs) to identify interpretable underperforming subsets of data and formulate hypotheses regarding the cause of observed performance disparities. We introduce a novel SDM and apply it in a case study on the classification of pneumothorax and atelectasis from chest x-rays. Our study demonstrates the effectiveness of SDMs in hypothesis formulation and yields an explanation of previously observed but unexplained performance disparities between male and female patients in widely used chest X-ray datasets and models. Our findings indicate shortcut learning in both classification tasks, through the presence of chest drains and ECG wires, respectively. Sex-based differences in the prevalence of these shortcut features appear to cause the observed classification performance gap, representing a previously underappreciated interaction between shortcut learning and model fairness analyses.

**Keywords:** Slice Discovery Methods · Algorithmic Fairness · Shortcut Learning · Chest X-ray · Model Debugging

## 1 Introduction

Machine learning models have shown great promise in medical image-based diagnosis, sometimes with performance claims that rival human experts. However, reported performance may overstate these models' clinical utility and safety [29]. Specifically, models may underperform or fail systematically on critical subsets of data even while overall average accuracy remains high. In computer vision research, such subsets are called *underperforming slices* or *blind spots* [7,22],
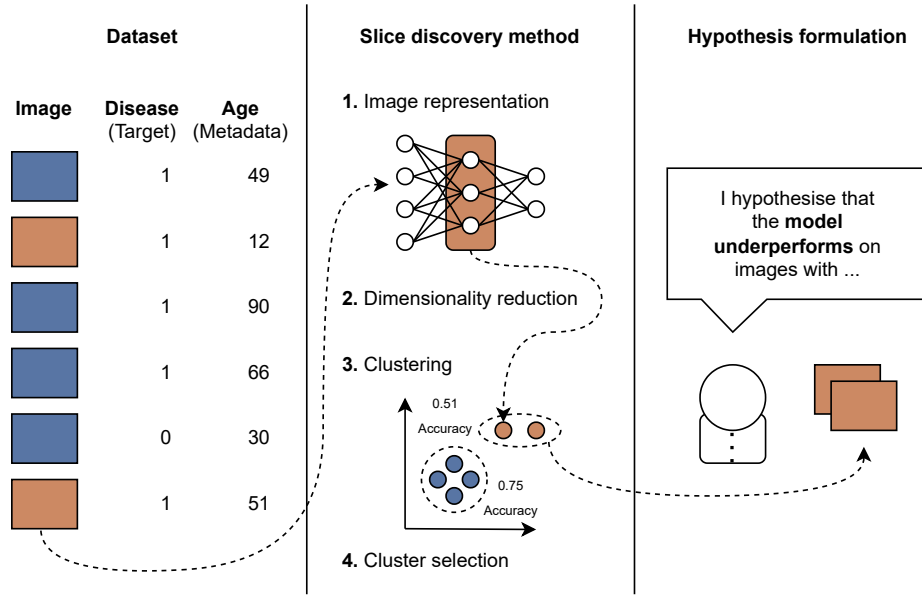
**Fig. 1.** A general overview of the key elements of slice discovery methods.

where 'slice' refers to a subset of samples with similar characteristics, such as an attribute familiar to a domain expert.[1]

In medical image analysis, a model might be underperforming on a slice of patients for a wide range of reasons, including group under-representation, increased input or label noise, fundamental differences in the difficulty of the prediction task, and shortcut learning [5,11,19,21]. Performance disparities between patient groups have been observed in many medical imaging domains [4,8,16,24,32], raising concerns about the potential unfairness resulting from the application of such models. However, properly *mitigating* such performance disparities requires identifying their root cause, which is often challenging [18,21,28]. The challenge is further compounded by the fact that the feature that causally distinguishes high-performing from low-performing patients is often unknown and, thus, not annotated. This renders simple subgroup analyses based on available metadata insufficient for identifying the causes of performance disparities.

To address the issue of unknown distinguishing features, various methods for the unsupervised discovery of underperforming slices have been proposed in the computer vision literature. Such methods are variously known as *Slice Discovery Methods* (SDMs) [7] or *Blindspot Discovery Methods* (BDMs) [22]. Typically, these methods perform a cluster analysis on the input space and then

---

[1] In the medical imaging literature, the term 'slice' commonly refers to a two-dimensional cross-section within three-dimensional volumetric data. We are adopting a differing terminology from earlier work on SDMs originating outside of the medical image analysis field. We apologize for the unfortunate clash of terminology.

select poorly performing clusters, or *slices* of data, for further analysis; refer to Fig. 1 for a high-level overview. SDMs can aid machine learning practitioners and domain experts in identifying underperforming sets of data, as well as in forming hypotheses about the *causes* of this underperformance. With few exceptions [19], SDMs have not yet found widespread use in the medical imaging domain.

In this study, we explore the use of SDMs for the analysis of performance disparities in medical image analysis. Our contributions are twofold. First, we provide a general overview of SDMs in medical image analysis and we propose a novel SDM, rigorously justifying all of our design choices. We demonstrate the effectiveness of our proposed SDM for hypothesis formulation in a case study of pneumothorax and atelectasis classification on two public chest X-ray datasets (NIH-CXR 14 [27] and CheXpert [9]). Second, by further analyzing the hypotheses generated using our SDM, we show that chest drain shortcut learning causes a previously unexplained yet variously reproduced performance gap between male and female subjects in pneumothorax classification. This constitutes an important link between shortcut learning and model fairness analyses that has, to the authors' knowledge, not been described before. In addition, using our SDM, we discover a new shortcut feature (the presence of ECG cables) that may explain male–female performance disparities in atelectasis classification.

## 2   Related work

### 2.1   Bias and shortcuts in chest x-ray analysis

Algorithmic fairness in medical image analysis, and performance disparities between patient groups in particular, have recently come under rapidly increasing scrutiny [4,14,19,23,24,28]. In this context, the fairness of chest x-ray-based disease classification models has received particularly broad attention [1,8,14,24,28,30]. Larrazabal et al. [14] demonstrated that such models had better classification performance for a particular patient group (based on biological sex) if that group was represented in higher proportions of the training dataset. While not the focus of their study, their results also indicated significant differences between model performance on male and female subjects, with the classification models performing better for either group in different diseases. These (sometimes large) performance gaps persisted even in the case of sex-balanced training sets. This observation prompted Weng et al. [28] to investigate the hypothesis that biological sex differences were causing these unexplained performance gaps. Based on their results, the authors dismissed breast shadows as a factor, but other biological sex differences contributing to performance gaps remain uncertain, leading to an unexplained gender disparity. Zhang et al. [30] employed standard algorithmic fairness mitigation approaches to the chest x-ray case, finding that simple group balancing was one of the most robust approaches – which did not, however, mitigate the performance gaps observed by Larrazabal et al. [14].

In a separate development, it has been widely demonstrated that chest x-ray-based disease classification models are prone to relying on shortcut learning [1,5,11,19]. Both Oakden-Rayner et al. [19] and Jiménez-Sánchez et al. [11]

**Table 1.** A summary of slice discovery methods. Clf: The representation used by the classification model under analysis. Adapted from Plumb et al. [22].

| Method | Rep. | Dim. reduction | Clustering |
|---|---|---|---|
| Algorithmic measurement [19] | Clf | | KNN |
| MultiAccuracy Boost [13] | VAE | | Rigid/decision-tree regression |
| GEORGE [25] | Clf/BiT emb. | UMAP (d=1,2) | GMM |
| Spotlight [6] | Clf | | Optimization problem |
| Planespot [22] | Clf | scvis (d=2) | GMM |
| Domino [7] | CLIP | PCA (d=128) | GMM |
| Failure mode distillation [10] | CLIP | | SVM |
| Bias-Aware Hierarchical Clustering [17] | Clf | UMAP(d=2) | Mod. K-Means |
| **Proposed SDM (Ours)** | Clf | FC layer (d=128) | GMM |

demonstrate how pneumothorax classification tends to rely heavily on the presence of chest drains, which represent the standard treatment for pneumothorax. Connecting the two challenges of shortcut learning and fairness, many authors have raised concerns about the potential for deep learning models to exploit spurious correlations between sensitive attributes, such as age, gender, or ethnicity, and the prediction target [1,8]. To the authors' knowledge, the fact that shortcut learning relying on *non*-sensitive features (such as the presence of chest drains) can explain performance disparities between sensitive groups (such as gender groups) has not been discussed explicitly before. Notably, Jiménez-Sánchez et al. [11] took important first steps in this direction, by differentially reporting the effect of shortcut learning on different gender groups.

## 2.2   Slice discovery methods

Slice discovery methods (SDMs) are a recently emerging tool for the performance analysis and subsequent improvement of deep learning models. In particular, they aim to solve the problem that the features that identify underperforming groups of inputs might not be known a priori. To this end, SDMs typically perform unsupervised clustering of the input data, in order to identify semantically similar 'slices' of data that the model under analysis performs poorly on. In more detail, SDMs usually consist of the following steps: (1) the input data is embedded into a latent representation space, (2) some SDMs perform dimensionality reduction, (3) an unsupervised learning method, such as clustering, is used to extract slices of the data, and (4) the extracted slices are prioritized based on a performance metric, such as accuracy. In table 1, adapted from Plumb et al. [22], we summarize previously proposed SDMs and their respective high-level design. The most common methods to extract slices are clustering algorithms, namely the Gaussian Mixture Model (GMM) and its variants. The most common choice of image embedding is the latent space representation computed by the classification model under scrutiny, i.e., its penultimate layer's outputs. However,

recent methods instead utilize multi-modal pre-trained models like CLIP to enable the generation of text descriptions for extracted slices [7,10]. Interestingly, the crucial dimensionality reduction step has received relatively little attention yet, as recently pointed out by Plumb et al. [22]. Possibly due to their relatively recent emergence, SDMs have not yet been widely applied in the medical image domain. In this regard, the work of Oakden-Rayner et al. [19] represents a very notable early exception that precedes more recent SDM developments.

## 3   Methodology

### 3.1   Proposed slice discovery method

This section introduces our proposed SDM and motivates our design choices. The proposed method consists of the following four steps, following Fig. 1:

**Image representation.** We use the image representation computed by the penultimate layer of the classification model under scrutiny. As opposed to SDMs that use a separate model to obtain image embeddings, our approach relies only on information available to the model's final classification layer.

**Dimensionality reduction.** We insert a single fully connected layer with d-dimensional output (and a sigmoid activation layer) between the classification model's penultimate layer and its final classification head. We train (just) this additional layer following the same procedure that was used for the classification model itself. Similarly to the previous step and contrary to standard choices such as PCA, t-SNE, or UMAP, our approach is *supervised* and preferably preserves information relevant to the model's predictions.

**Clustering.** We use a Gaussian Mixture Model (GMM) for clustering and the Bayesian Information Criterion (BIC) for choosing the number of clusters. We cluster disease-positive and -negative samples separately to extract clusters of the same error type, similar to Oakden-Rayner et al. [19].

**Cluster selection.** We propose using the Brier score (BS), a proper scoring rule, to prioritize under- and overperforming slices, equivalent to mean squared error between model confidence and binary target labels. The main motivation for our choice is that, as opposed to classification accuracy, the BS is threshold-independent. In addition, it captures both the model's discriminative ability and its calibration [2]. As opposed to AUROC [12], per-cluster BS can be meaningfully compared between groups, and as opposed to many calibration metrics [20,23], it can be meaningfully compared between clusters of different sizes. We quantify BS uncertainty by simple bootstrapping of each cluster and select the cluster with the lowest 97.5-quantile as the best, and the cluster with the highest 2.5-quantile as the worst.

### 3.2   Datasets

We consider a case study on two public datasets, NIH-CXR14 [27] and CheX-pert [9]. Both datasets slightly over-represent male subjects. We reused chest

drain labels previously crowd-sourced from both radiologists and nonexperts [3,19,11], including 3543 random cases with pneumothorax in the NIH dataset and 972 cases with and without pneumothorax in CheXpert. For NIH, we observe a larger prevalence of chest drains among pneumothorax-positive male subjects compared to female subjects (49.5% vs. 42.8%). For CheXpert, we observe a larger prevalence of chest drains among pneumothorax-negative male subjects (23.0% vs. 14.7% in females) but comparable chest drain prevalence across sexes for pneumothorax-positive subjects (50.5% in males vs. 50.2% in females).

### 3.3   Experiments

We conduct a case study on the pneumothorax classification task, following the experimental setup of Weng et al. [28]. Specifically, to reduce the potential for label noise to affect our analyses, we select one sample per patient, with a preference for pneumothorax-positive samples and an equal sex ratio. Withholding the chest drain-annotated samples, we split the datasets into 60%/10%/30% train, validation, and test sets, resampling the splits ten times. We resize the images to 224x224 pixels and train a ResNet50 (Adam optimizer, learning rate $10^{-6}$, batch size 64, 20 epochs). We use data augmentation for the training dataset, including horizontal flipping, rotation up to 15 degrees, and scaling up to 10% with a 50% probability for each augmentation. We then carry out our proposed SDM and report the distribution of comorbidities and chest drains (for annotated chest drain samples).[2] We repeat the same analyses for atelectasis classification. Based on our findings, we conduct further post-hoc analyses in both cases. Our source code is publicly available at https://github.com/volesen/slicing-through-bias.

## 4   Results

*Pneumothorax classification.* Fig. 2 and fig. 4 show the best- and worst-performing slices based on the Bier score in the NIH-CXR14 and CheXpert datasets. In both datasets, the underperforming slices for pneumothorax-negative samples have a lower-than-average chest drain proportion, while the opposite holds for pneumothorax-positive cases.

Based on these results (and based on prior work [11,19]), we hypothesized that the presence of chest drains is used to classify pneumothorax. Indeed, we observe that computed model confidence (the softmax output of the model for the disease-positive class) is consistently higher in samples with chest drains across datasets, sexes, and pneumothorax labels (positive/negative), lending strong support to this hypothesis (fig. 5). Furthermore, as both datasets have varying chest drain prevalences by sex, we hypothesized that this could contribute to

---

[2] For the dimensionality reduction, we used d=128 following similar previous work [7]. Results with d=10 were comparable. The gap statistic [26] indicates that clustering does indeed occur in the reduced space. The effect of the additional dimensionality reduction layer on the model's classification performance was negligible in terms of test accuracy and AUROC.
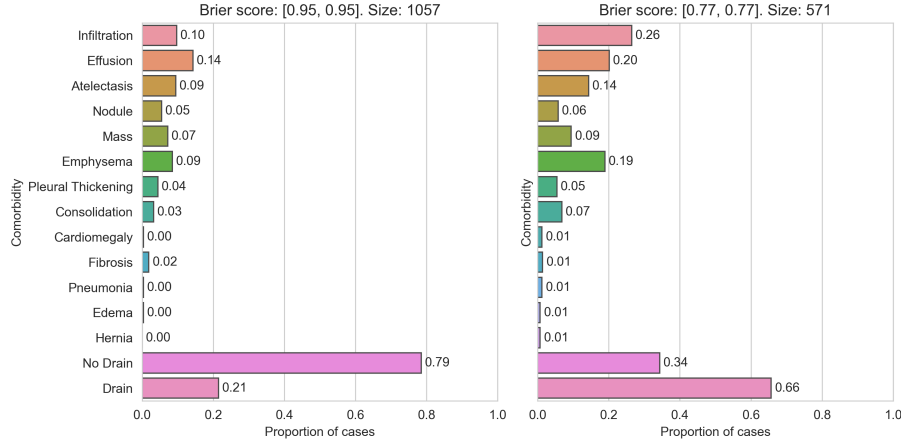
Brier score: [0.95, 0.95]. Size: 1057          Brier score: [0.77, 0.77]. Size: 571

**Fig. 2.** The comorbidity and chest drain distribution in pneumothorax-positive chest drain annotated samples of NIH-CXR14 for the worst-performing (left column) and best-performing (right column) slices by Brier score. The pneumothorax-negative case is omitted as chest drain annotations were not available for these samples in NIH-CXR14.

the gender performance gap. To assess this hypothesis, we determined AUROC by sex, firstly in a test set following the natural distribution of chest drains, and secondly a test set with equalized chest drain prevalence in the male and female populations (Fig. 3). We observe a significant difference in performance in the first case ($p < 0.001$) but no significant difference in the second, chest-drain-balanced case ($p > 0.1$), indicating that chest drain shortcut learning is the cause of a large part of the male–female performance gap observed in prior work [14,28]. Statistical significance was assessed using a Mann-Whitney U-test. It should be noted that the differences in comorbidities noticeable in Fig. 2 represent another potential explanation for the observed performance disparities. While we did not further explore this hypothesis here, our additional experiments related to the chest drain shortcut hypothesis (discussed above) seem to indicate that this indeed explains the bulk of the observed performance gaps.

*Atelectasis classification.* We repeated our SDM analysis for atelectasis. (Slice statistics not shown here due to space constraints.) Prompted by a visual inspection of the recordings in the best and worst-performing slices, we hypothesized that the two slices differed substantially in their prevalence of ECG cables in the recordings. To test this hypothesis, we randomly selected and (as non-experts) labeled 100 samples each for the best- and worst-performing slices on atelectasis-positive and -negative cases according to whether they displayed ECG cables or not. Of atelectasis-negative samples, 95% of the labeled recordings in the worst-performing slice contained ECG cables, compared to 10% in the best-performing
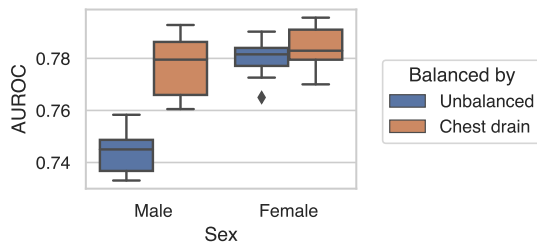
**Fig. 3.** AUROC on CheXpert with male and female test subjects on pneumothorax prediction, following the natural ('unbalanced') distribution of chest drains and balanced by chest drain presence across ten samplings of the train-validation-test sets.

slice. Of the atelectasis-positive cases, 50% of the labeled recordings in the worst-performing slice had ECG cables, compared to 99% in the best-performing slice.

## 5    Discussion and conclusion

We have proposed a novel slice discovery method (SDM), which differs from previously proposed methods in several key elements. Both in the representation extraction as well as in the (supervised) dimensionality reduction step, we prioritize only using information available to the classification model under test. In addition, we propose using the Brier score (BS) for selecting highly and poorly performing clusters because it is threshold-independent, accounts for both discriminative ability and calibration, and enables meaningful comparisons between clusters of different sizes. In a case study on chest x-ray-based disease classification, our SDM successfully recovered a previously known case of shortcut learning (chest drains for pneumothorax classification) and suggested a new, previously unknown case (ECG cables for atelectasis classification). The latter case also demonstrated another benefit of using SDMs: reducing the required labeling efforts, because these can be specifically targeted at the worst and best clusters. Our case study shows that our proposed SDM, and SDMs in general, can aid researchers in generating hypotheses regarding the causes of model underperformance on subsets of data, which is crucial for leveling *up* performance [21].

Consistent with the observations made by Weng et al. [28], our findings challenge the notion that biological differences are the primary driver of the previously observed [14] but unexplained male–female performance gaps in chest x-ray-based disease classification. Instead, our results suggest that shortcut learning in conjunction with a difference in chest drain prevalence between males and females causes the observed performance disparity. This newly gained knowledge opens up the possibility for the targeted application of shortcut learning mitigation techniques, instead of relying on blind group performance equalization approaches that often result in leveling *down* performance [21,30,31].

While our study builds upon the chest drain annotations of [3] and [11], for which the former show a high level of agreement between expert and non-expert

annotations, caution is warranted. Non-expert annotations, although showing agreement, may not necessarily represent the ground truth or offer a representative sample (images without consensus between multiple labelers were disregarded [11]). This caution is emphasized by a significant difference in the prevalence of chest drains among pneumothorax-positive samples in NIH-CXR14 between studies [3,19]. Moreover, the mitigation of shortcut learning is a highly active research area, and mitigating many different shortcuts simultaneously, and with limited label availability, remains challenging [15]. Finally, the interpretation of identified slices in more challenging cases represents a crucial challenge for future research. Both chest drains and ECG cables are visible to the human eye (though maybe not the non-expert), but other potentially problematic features may not be.

## Acknowledgements

## References

1. Brown, A., Tomasev, N., Freyberg, J., Liu, Y., Karthikesalingam, A., Schrouff, J.: Detecting shortcut learning for fair medical AI using shortcut testing. Nature Communications **14**(1) (2023)
2. Bröcker, J.: Reliability, sufficiency, and the decomposition of proper scores. Quarterly Journal of the Royal Meteorological Society **135**(643), 1512–1519 (2009)
3. Damgaard, C., Eriksen, T.N., Juodelyte, D., Cheplygina, V., Jiménez-Sánchez, A.: Augmenting chest x-ray datasets with non-expert annotations
4. Daneshjou, R., Vodrahalli, K., Novoa, R.A., Jenkins, M., Liang, W., Rotemberg, V., et al.: Disparities in dermatology AI performance on a diverse, curated clinical image set. Science Advances **8**(32) (2022)
5. 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)
6. d'Eon, G., d'Eon, J., Wright, J.R., Leyton-Brown, K.: The Spotlight: A general method for discovering systematic errors in deep learning models. In: ACM FAccT. p. 1962–1981 (2022)
7. Eyuboglu, S., Varma, M., Saab, K.K., Delbrouck, J.B., Lee-Messer, C., Dunnmon, J., et al.: Domino: Discovering systematic errors with cross-modal embeddings. In: ICLR (2022)
8. Glocker, B., Jones, C., Roschewitz, M., Winzeck, S.: Risk of bias in chest radiography deep learning foundation models. Radiology: Artificial Intelligence **5**(6) (2023)
9. Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., et al.: CheXpert: a large chest radiograph dataset with uncertainty labels and expert comparison. In: AAAI/IAAI/EAAI (2019)

10. Jain, S., Lawrence, H., Moitra, A., Madry, A.: Distilling model failures as directions in latent space. In: ICLR (2023)
11. Jiménez-Sánchez, A., Juodelyte, D., Chamberlain, B., Cheplygina, V.: Detecting Shortcuts in Medical Images - A Case Study in Chest X-Rays. In: ISBI. IEEE, Cartagena, Colombia (2023)
12. Kallus, N., Zhou, A.: The fairness of risk scores beyond classification: Bipartite ranking and the xAUC metric. In: NeurIPS. vol. 32. Curran Associates, Inc. (2019)
13. Kim, M.P., Ghorbani, A., Zou, J.: Multiaccuracy: Black-box post-processing for fairness in classification. In: AIES. pp. 247–254. ACM (2019)
14. 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)
15. Li, Z., Evtimov, I., Gordo, A., Hazirbas, C., Hassner, T., Ferrer, C.C., et al.: A whac-a-mole dilemma: Shortcuts come in multiples where mitigating one amplifies others. In: CVPR. pp. 20071–20082 (2023)
16. Lin, M., Li, T., Yang, Y., Holste, G., Ding, Y., Van Tassel, S.H., et al.: Improving model fairness in image-based computer-aided diagnosis. Nature Communications **14**(1) (2023)
17. Misztal-Radecka, J., Indurkhya, B.: Bias-aware hierarchical clustering for detecting the discriminated groups of users in recommendation systems. Information Processing & Management **58**(3), 102519 (2021)
18. Mukherjee, P., Shen, T.C., Liu, J., Mathai, T., Shafaat, O., Summers, R.M.: Confounding factors need to be accounted for in assessing bias by machine learning algorithms. Nature Medicine **28**(6), 1159–1160 (2022)
19. Oakden-Rayner, L., Dunnmon, J., Carneiro, G., Re, C.: Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In: CHIL. pp. 151–159. ACM (2020)
20. Petersen, E., Ganz, M., Holm, S.H., Feragen, A.: On (assessing) the fairness of risk score models. In: FAccT. ACM (2023)
21. Petersen, E., Holm, S., Ganz, M., Feragen, A.: The path toward equal performance in medical machine learning. Patterns **4**(7) (2023)
22. Plumb, G., Johnson, N., Cabrera, A., Talwalkar, A.: Towards a more rigorous science of blindspot discovery in image classification models. Transactions on Machine Learning Research (2023)
23. Ricci Lara, M.A., Mosquera, C., Ferrante, E., Echeveste, R.: Towards Unraveling Calibration Biases in Medical Image Analysis, pp. 132–141. Springer Nature Switzerland (2023)
24. Seyyed-Kalantari, L., Zhang, H., McDermott, M.B.A., Chen, I.Y., Ghassemi, M.: Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. Nature Medicine **27**(12), 2176–2182 (2021)
25. Sohoni, N., Dunnmon, J., Angus, G., Gu, A., Ré, C.: No subclass left behind: Fine-grained robustness in coarse-grained classification problems. In: NeurIPS. vol. 33, pp. 19339–19352 (2020)
26. Tibshirani, R., Walther, G., Hastie, T.: Estimating the number of clusters in a data set via the gap statistic. Journal of the Royal Statistical Society Series B: Statistical Methodology **63**(2), 411–423 (Jul 2001). https://doi.org/10.1111/1467-9868.00293
27. 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 classification and localization of common thorax diseases. In: CVPR. pp. 2097–2106 (2017)

28. Weng, N., Bigdeli, S., Petersen, E., Feragen, A.: Are sex-based physiological differences the cause of gender bias for chest x-ray diagnosis? In: MICCAI Workshop on Fairness of AI in Medical Imaging. pp. 142–152. Springer (2023)
29. Wynants, L., Van Calster, B., Collins, G.S., Riley, R.D., Heinze, G., Schuit, E., et al.: Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. BMJ p. m1328 (2020)
30. Zhang, H., Dullerud, N., Roth, K., Oakden-Rayner, L., Pfohl, S., Ghassemi, M.: Improving the fairness of chest x-ray classifiers. In: Conference on Health, Inference, and Learning. pp. 204–233. PMLR (2022)
31. Zietlow, D., Lohaus, M., Balakrishnan, G., Kleindessner, M., Locatello, F., Schölkopf, B., Russell, C.: Leveling down in computer vision: Pareto inefficiencies in fair deep classifiers. In: CVPR. IEEE (2022)
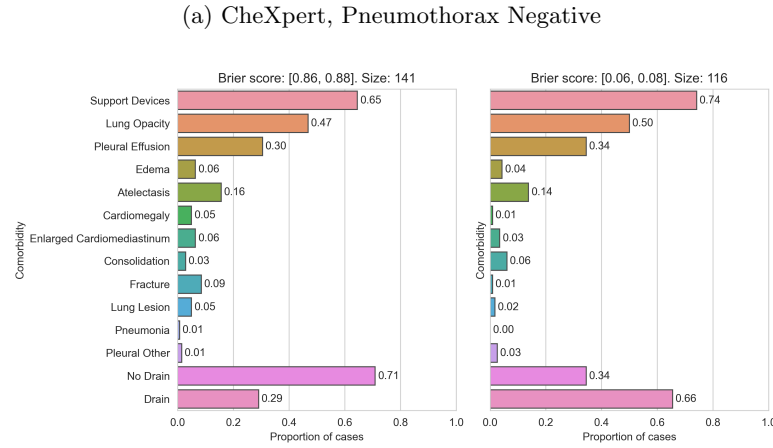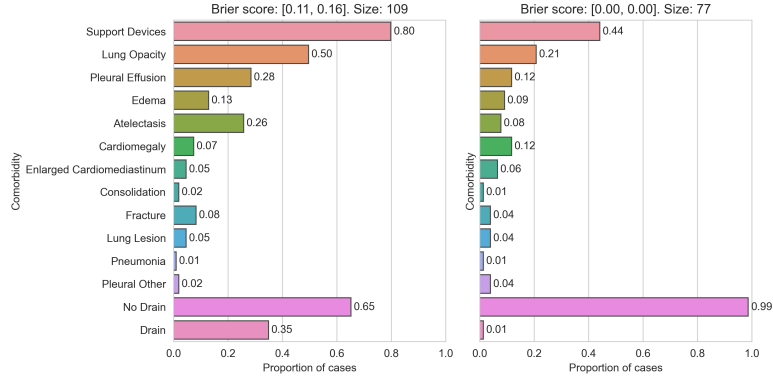32. Zong, Y., Yang, Y., Hospedales, T.: MEDFAIR: Benchmarking fairness for medical imaging. In: ICLR (2023)

(a) CheXpert, Pneumothorax Negative



(b) CheXpert, Pneumothorax Positive

**Fig. 4.** The comorbidity and chest drain distribution in (a) pneumothorax-negative (top row) and (b) pneumothorax-positive chest drain annotated samples of CheXpert for the worst-performing (left column) and best-performing (right column) slices by upper 95% bootstrapped Brier scores.
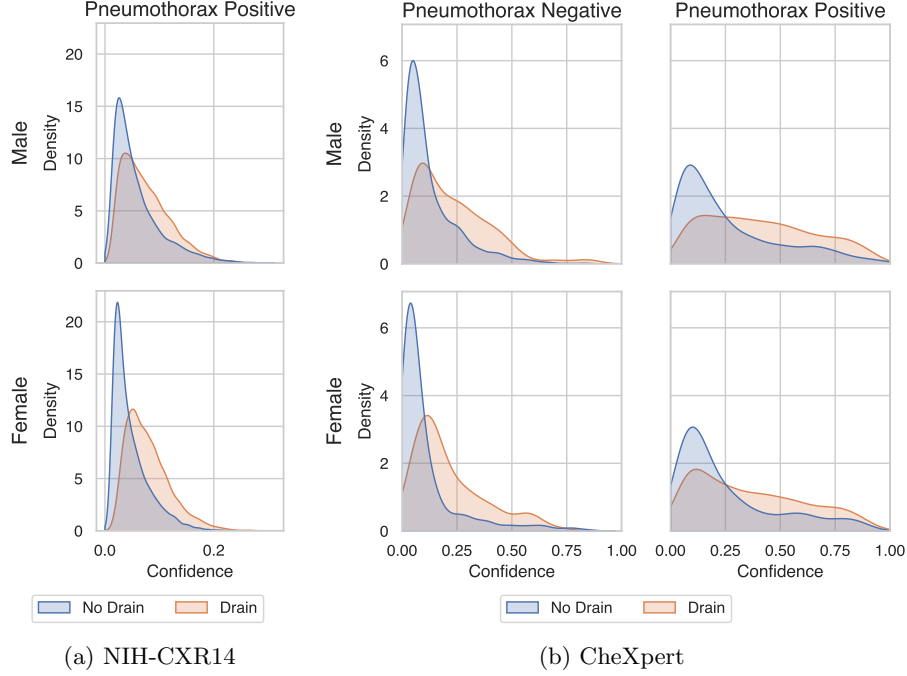
(a) NIH-CXR14

(b) CheXpert

**Fig. 5.** Distribution of confidences (the softmax output of the model for the disease-positive class) for pneumothorax classification by sex, presence of pneumothorax, and chest drain for NIH-CXR14 (left) and CheXpert (right). Throughout, subjects without chest drains are more likely to be classified as pneumothorax-negative.
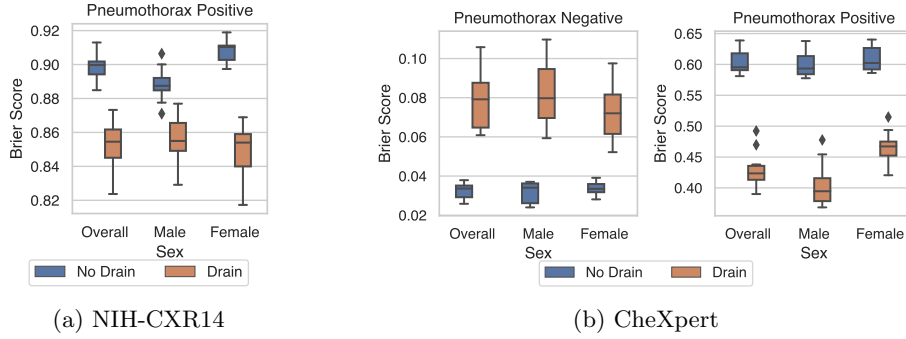


(a) NIH-CXR14

(b) CheXpert

**Fig. 6.** Brier scores for pneumothorax classification for NIH-CXR14 (left) and CheXpert (right) across ten samplings of the train-validation-test sets. Throughout, non-pneumothorax patients with chest drains and pneumothorax patients without chest drains are underperforming compared to pneumothorax-positive subjects with chest drains and pneumothorax-negative subjects without chest drains.