

Hidden Stratification in Medical Imaging: Bias and Subgroup Effects

1. Introduction

The paper “Hidden Stratification Causes Clinically Meaningful Failures in Machine Learning for Medical Imaging” identifies major clinical failings in deep learning models for medical image classification due to their exhibiting an overall accuracy that is misleadingly high and masking the poor performance on clinically relevant subgroups by the authors. The hidden stratification concept introduced by the authors is said to occur when unappreciated subsets of data within a class have been innately un-set up to perform quite differently from others in terms of properties that could adversely affect model generalization. This has been demonstrated on three datasets: (1) a pelvic X-ray dataset from the Royal Adelaide Hospital for hip fracture detection (2,000+ images); (2) the MURA dataset for classification of abnormalities on musculoskeletal X-rays (40,000+); and (3) the CXR14 dataset for the identification of thoracic pathologies focusing on pneumothorax (112,000+ chest X-rays). The models used are DenseNet architectures (e.g., DenseNet-121 and DenseNet-169) trained for binary or multi-label classification. The application is clinical diagnostics, where such failures could prioritize common, benign cases over rare, life-threatening cases, stressing the need for subgroup-aware evaluations in healthcare AI to ensure safety and equality.

2. Source of Bias

Hidden stratifications arise mainly from incomplete schema design during the labeling of data. Broad categories thus become a collection of diverse subsets that escape explicit annotation. This occurs in the paper during data collection and preprocessing stages. Continuing with the CXR14 dataset, the pneumothorax labels lump together treated (with chest drains) and untreated cases, resulting in spurious correlations—chest drains are non-causal artifacts from treatment, but models learn these as discriminant features. This represents aggregation bias, whereby heterogeneous groups are lumped together, masking differences among subgroups.

Sampling biases arise in the hip fracture dataset due to under-sampling: low-prevalence subsets such as cervical fractures (13% of total cases) or subtle fractures (6%) can hardly be represented, leading to representation bias. These biases are also accepted in the paper as samples are collected from hospital records: sampling is skewed naturally towards more commonly seen variants (like peritrochanteric fractures, statistically at 50%) since all other serious diseases are rare. This is worsened by preprocessing if stratification of images is not performed so that models learn to fit to the dominating pattern.

The MURA dataset highlights the labelling bias of measurement “abnormal” fractures (92% reporting sensitivity), hardware (85%), and degenerative joint disease (DJD, little downward sensitivity at 60%). Clinical reports often label only the grossest of fractures and in doing so ignore subtle or irrelevant findings and thereby downplay the diagnosis for DJD, whose consequences are vague and non-urgent. This agrees with concepts in the course such as label bias, where noisy or incomplete annotations unduly affect subpopulations with subtle features.

In summary, these biases stem from the realities of the real world, where clinical priorities favor acute cases, giving rise to hidden subgroups formed through incomplete representation and measurement errors.

3. Impact on Metrics

Hidden stratification leads to distortion of evaluation metrics, permitting dominant subgroups to inflate general statistics, thereby hiding a poor performance of minorities. The ROC curve area and two aggregate measures—sensitivity (recall) and positive predictive value (PPV)—have been used in the paper, although they are sensitive to class imbalances. Overall sensitivity in the hip fracture dataset is 0.981 (AUC 0.994), while diagnostic sensitivity drops to 0.911 with respect to cervical fractures and 0.900 for the subtle ones. The problem in segmenting these subgroups is that they have a low prevalence and rely on subtle discriminative features—cervical fractures do not exhibit any obvious displacement

and tend to get lost in the noise. A sizeable subgroup like the pertrochanteric is falsely assumed to be accurate, based on a sensitivity of 0.997—and that really throws the majority into thinking they have expert-level accuracy. The above just reiterates, once again, why accuracy may be a misleading definition in an imbalanced dataset: It cares for bulk rather than clinical severity.

For the MURA study, abnormality detection yielded an overall AUC of 0.91, with hardware (easy-to-spot metallic implants) scoring 0.98 and DJD lagging behind at 0.76. DJD underperforms due to poor label quality (sensitivity 0.60) and subtle features like joint narrowing, which the models confuse with normals. Fractures perform slightly better (AUC 0.86) because these fractures have clear breaks and high label accuracy (0.92) in the DJD subgroup. The subgroups have high prevalence (DJD at 43%) but low discriminability and drag metrics, but overall AUC hides this by averaging across easier cases. This shows how ROC AUC, though threshold-independent, can mask an uneven playing field when subgroups vary dramatically in feature subtlety and label noise—models shine on the "simply" easy hardware (11% prevalence) but flop on a clinically profound DJD.

The CXR14 example of pneumothorax is dramatic: overall AUC 0.87, but 0.94 with chest drains versus 0.77 otherwise. When drains (80% prevalence) exist, they act as spurious correlates, yielding performance gains via artifacts (visible tubes as proxies), with PPVs increased by 30% (0.90 vs. 0.60). In the untreated cases—subtler ones, without artifacts—the models' reliance on coincidental non-causal cues will spoil further. The asymmetry between the subgroups becomes very crucial: An untreated pneumothorax is very lethal, yet it is dominated by comforting metrics of treated (benign) cases. In theory, because great data brings in the hidden correlations as per Selbst (2017), it makes aggregate metrics unsuitable for safety-critical applications. In healthcare, this could have regulatory effects, where models appear at the level of the radiologist while doing real harm to vulnerable patients.

4. Mitigation & Creativity

The three mitigation strategies recommended by Oakden-Rayner et al. include schema completion, error auditing, and algorithmic measurement.

- Schema completion assigns explicit subgroup labels to a testing corpus (e.g., fracture types or presence/absence of drains) such that ASER could report on sub performance across subgroups, though it is limited by available human resources as well as incomplete medical taxonomies.
- An error audit systematically reviews misclassified samples to look for patterns or recurrent failure modes; such an approach would help discover hidden strata post hoc.
- Algorithmic measurements use unsupervised clustering on learned feature embeddings to potentially identify subgroups with different error rates. Reducing human workload—again, the clinical usefulness requires some validation.

The new enhancement would then be a hybrid bias detection pipeline that uses explainable AI along with active learning. More specifically, saliency maps or Grad-CAM would identify regions relied upon by the model when spurious features like chest drains dominate activations, and the system flags the image for human review. That feedback loop could provide guidance in targeted relabeling or model retraining.

Along with this, subgroup-aware measures would also need to be incorporated in evaluation metrics—that is, performance metrics such as worst-case AUC or equal opportunity difference—ensuring that the measurement upholds fairness also in hidden strata. This means that no subgroup, instead of the average subgroup, would additionally be required to meet an acceptable standard in performance; this would make deployment safer.

At last, synthetic data augmentation through generative models would balance low-prevalence subgroups, rendering the model even less sensitive to dominant patterns. Well articulating these strategies into a proactive approach will enable the detection, quantification, and sabotage of hidden stratification bias in medical imaging systems.