

---

# On the Impact of Gender Bias in Medical Imaging Classifiers for Computer-aided Diagnosis

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Research Problem

Artificial intelligence (AI) influences almost every aspect of our daily life. In particular, the rise of AI in healthcare during the last few years is changing the way medical doctors diagnose, especially when dealing with medical images. AI systems can not only augment the information provided by such images with useful annotations [1], but they are also taking autonomous decisions by performing computer assisted diagnosis (CAD) [2, 3]. In this context, care must be taken in the risk associated with error and misbehaviour of algorithms, especially since their decision comprises a delicate domain such as health care.

Although the interest in performing fair and unbiased evaluations of AI medical systems exists since the 80's [4], the ethical and moral aspects of AI have gained relevance in the last few years, showing that human bias, such as gender and race bias may be inherited by AI systems in multiple contexts [5, 6]. In the last years, the research community of gendered innovations [7] has worked to create awareness and integrate sex and gender analysis into all phases of basic and applied research. However, to date, there is no study reflecting such analysis in the context of medical imaging and computer assisted diagnosis.

In this work, we perform the first large-scale study that quantifies the influence of gender imbalance in medical imaging datasets used to train AI-based CAD systems. We employ a classification model based on deep neural networks, which achieves state-of-the-art results when diagnosing 14 common thoracic diseases using X-ray images [8]. We analyze the performance over male and female patients when the model is trained with different gender unbalance ratios, providing empirical evidence about the bias acquired by such systems.

## 2 Experiments

We use a frontal view chest X-ray image dataset, that contains over 100.000 images with 14 different common pathologies [8]. The dataset has 63340 images of male patients and 48780 female patients. All images were automatically labeled using natural language processing tools to extract such information from radiology reports (see [9] for a detailed method description).

Given a frontal X-ray image, we trained CAD models to predict the presence or absence of the 14 thoracic diseases considering male-only and female-only training datasets. We then evaluated both classifiers in male and female patients separately, and reported their performance using the well-known area under the receiving operating curve (AUC) [10]. Since we are focusing on the effect of gender imbalance, we guaranteed by construction that male and female folds include the same number of pathological cases per class to avoid other sources of bias.

A Densely Connected Convolutional Neural Network (DenseNet) [11] architecture with 14 outputs representing the probability of each disease was used for classification. We adopted a Keras implementation of the DenseNet-121 (publicly available at: <https://github.com/brucechou1983/CheXNet-Keras>)

36 which has shown to achieve state-of-the-art results in X-ray image classification [8]. The network  
 37 has 121 convolutional layers and a final fully connected layer producing a 14-dimensional output,  
 38 after which we apply an element-wise sigmoid non-linearity. A model pre-trained on ImageNet [12]  
 39 was used to initialize the network. We trained it end-to-end using Adam optimizer with standard  
 40 parameters ( $\beta_1=0.9$  and  $\beta_2=0.999$ ), a batch size of 32 and an initial learning rate of 0.001 that was  
 41 decayed by a factor of 10 each time the validation loss plateaus after an epoch.

42 For every experiment, 20 models were trained using random sampling and the area under the curve  
 43 (AUC) was used to measure classification performance. For each random subset, male and female  
 44 patients were evaluated separately, considering for each gender 20% test, 70% training and 10%  
 45 validation, ensuring that images from one patient were not overlapping in different splits. In all cases,  
 46 we tested each model with a female test set and male test set independently.

### 47 3 Results and discussion

48 The experimental results shown in Figure 1 show a consistent decrease in performance when using  
 49 male patients for training and female for testing, and viceversa. In particular, we found that 16 out  
 50 of 28 cases (14 diseases per gender) present a significant decrease in performance, according to a  
 51 Mann-Whitney U test, considering a p-value  $< 0.01$ .

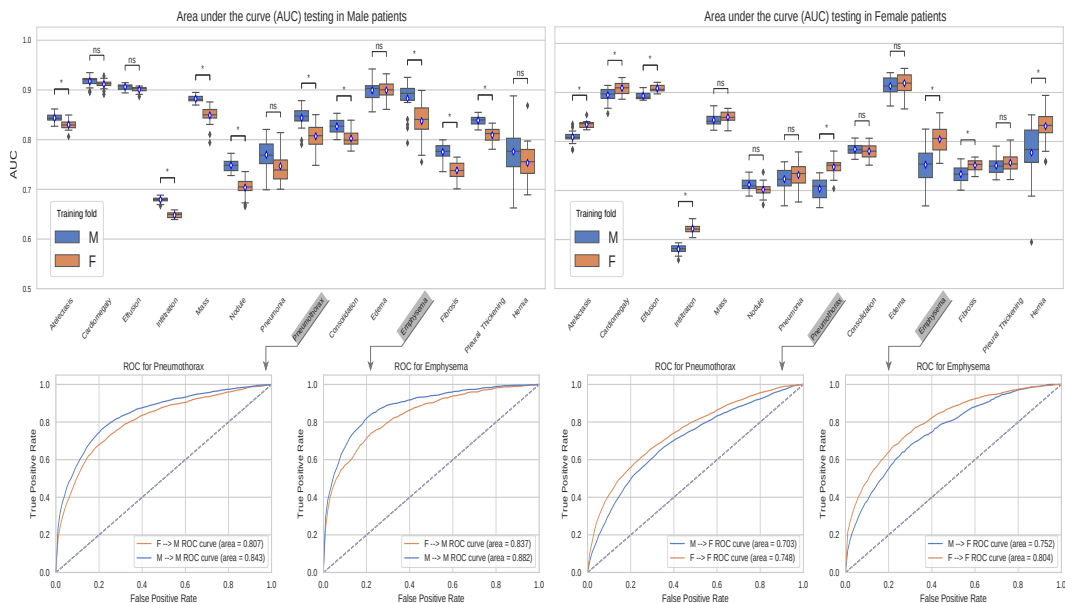


Figure 1: Top: AUC for each disease corresponding to a CNN trained only with male patients (blue) and trained only with female patients (orange). Bottom: ROC curve from specific diseases (p-value  $< 0.01$  according to a Mann-Whitney U test).

### 52 4 Conclusions

53 We have shown how algorithms can present gender bias if they are trained in unbalanced data sets,  
 54 producing serious misbehaviour for the under-represented class. This raises the alarm for national  
 55 agencies in charge of regulating and approving CAD systems, which should include explicit gender  
 56 balance and diversity recommendations. We also establish a new open problem for the academic  
 57 medical image computing community which needs to be addressed by novel algorithms endowed  
 58 with robustness to gender imbalance.

59 In the future we plan to develop algorithmic solutions to tackle gender bias problems in cases where  
 60 it is difficult to obtain balanced datasets with annotations. We plan to employ adversarial domain  
 61 adaptation techniques [13] to train classification models which are invariant to patient gender.

## 62 References

- 63 [1] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco  
64 Ciompi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I  
65 Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*,  
66 42:60–88, 2017.
- 67 [2] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and  
68 Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks.  
69 *Nature*, 542(7639):115, 2017.
- 70 [3] Jeffrey De Fauw, Joseph R Ledsam, Bernardino Romera-Paredes, Stanislav Nikolov, Nenad  
71 Tomasev, Sam Blackwell, Harry Askham, Xavier Glorot, Brendan O’Donoghue, Daniel Visentin,  
72 et al. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature*  
73 *medicine*, 24(9):1342, 2018.
- 74 [4] B Chandrasekaran. On evaluating artificial intelligence systems for medical diagnosis. *AI*  
75 *magazine*, 4(2):34–34, 1983.
- 76 [5] James Zou and Londa Schiebinger. Ai can be sexist and racist—it’s time to make it fair, 2018.
- 77 [6] Matthew Hutson et al. Even artificial intelligence can acquire biases against race and gender.  
78 *Science Magazine*, 10, 2017.
- 79 [7] Londa Schiebinger and Martina Schraudner. Interdisciplinary approaches to achieving gen-  
80 dered innovations in science, medicine, and engineering1. *Interdisciplinary Science Reviews*,  
81 36(2):154–167, 2011.
- 82 [8] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy  
83 Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, et al. Chexnet: Radiologist-level  
84 pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*,  
85 2017.
- 86 [9] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M  
87 Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-  
88 supervised classification and localization of common thorax diseases. In *Proceedings of the*  
89 *IEEE conference on computer vision and pattern recognition*, pages 2097–2106, 2017.
- 90 [10] Tom Fawcett. An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874, 2006.
- 91 [11] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected  
92 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern*  
93 *recognition*, pages 4700–4708, 2017.
- 94 [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
95 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*  
96 *recognition*, pages 248–255. Ieee, 2009.
- 97 [13] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François  
98 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural  
99 networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.