

# Disparate Model Performance and Stability in Machine Learning Clinical Support for Diabetes and Heart Diseases

Ioannis Bilonis, MSc<sup>1,2</sup>, Ricardo C. Berrios<sup>1</sup>, Luis Fernandez-Luque, PhD<sup>1</sup>,  
Carlos Castillo, PhD<sup>2,3</sup>

<sup>1</sup>Adhera Health, Santa Cruz, USA; <sup>2</sup>Universitat Pompeu Fabra, Barcelona, Spain;  
<sup>3</sup>ICREA, Catalonia, Spain

## Abstract

*Machine Learning (ML) algorithms are vital for supporting clinical decision-making in biomedical informatics. However, their predictive performance can vary across demographic groups, often due to the underrepresentation of historically marginalized populations in training datasets. The investigation reveals widespread sex- and age-related inequities in chronic disease datasets and their derived ML models. Thus, a novel analytical framework is introduced, combining systematic arbitrariness with traditional metrics like accuracy and data complexity. The analysis of data from over 25,000 individuals with chronic diseases revealed mild sex-related disparities, favoring predictive accuracy for males, and significant age-related differences, with better accuracy for younger patients. Notably, older patients showed inconsistent predictive accuracy across seven datasets, linked to higher data complexity and lower model performance. This highlights that representativeness in training data alone does not guarantee equitable outcomes, and model arbitrariness must be addressed before deploying models in clinical settings.*

## Introduction

In the realm of biomedicine, Artificial Intelligence (AI) methodologies, particularly Machine Learning (ML) models, are used as clinical support tools to systematically discern patterns and interdependencies among factors and outcomes within large datasets. ML has the potential to enhance healthcare provision by complementing, rather than supplanting, clinical judgment. It has demonstrated efficacy in the detection of skin cancer and diabetic retinopathy, among many other medical conditions [1, 2]. A paramount objective when deploying ML models is the assurance of health equity [3, 4]; thus, researchers and practitioners typically aim at attaining uniform model efficacy across diverse patient demographics [5]. The academic literature recommends an array of analytical tools for detecting biases, e.g., determining statistical dependencies between model outcomes, model errors, and specific subgroups [6], particularly those experiencing both historical and ongoing discrimination. Disparities detected in model performance are frequently attributed to deficiencies within the training datasets, typically lack of sufficient samples from those groups [7].

Algorithmic fairness, as a research field, studies how and to which extent algorithmic decision support systems can be free from discriminatory biases [8, 9]. Discrimination, in this context, means systematic disadvantages affecting socially salient groups [10]. These disadvantages arise from a complex combination of design choices made at different points in the construction of an ML processing pipeline. Discriminatory biases have been documented in basically all applications of ML and AI [11], including recruitment [12], machine translation [13] and face recognition [14], just to name a few.

In healthcare applications, prior research has identified algorithmic bias as a factor contributing to health disparities, highlighting the need for including Social Determinants of Health (SDoH) in ML to achieve health equity [15, 16]. For instance, in computer vision applications for medical imaging, biased data has been found to be a source of disparities in algorithmic outcomes [17, 18]. Differences in mortality prediction and X-ray diagnosis have been identified across racial/ethnic groups [19, 20], including discrepancies in burn identification and diabetic retinopathy identification in dark-skinned versus lighter-skinned patients [21, 22], and in an opioid misuse classifier, with more errors (false negatives) for dark-skinned patients [23]. In other cases, ML algorithms have predicted similar risk scores in both light- and dark-skinned patients, even though the dark-skinned patients had higher risk [24, 25]. There are many other examples, as this is an active research topic that to some extent is in its early stages [26, 27, 28, 29].

An in-depth knowledge of an ML application and of its context should inform this analysis [30]. In healthcare, the generalizability of AI algorithms across subgroups is critically dependent on training datasets, including factors such as representativeness, missing data, and outliers [31]. This suggests that some biases can be traced to datasets that underrepresent certain populations; using these unbalanced datasets as training data yields algorithmic models that

exhibit systematically unbalanced errors [32]. In this context, the augmentation of the dataset with additional samples from the underrepresented group, which frequently corresponds to groups that are socioeconomically disadvantaged or medically underserved, has been empirically demonstrated to mitigate discrepancy in model accuracy. This is the case of the seminal “Gender Shades” study [7, 33]. Similar results have been observed in the training set of a popular face detection benchmark dataset [34].

Differences in algorithmic performance are not always due to lack of representativeness. Signs and symptoms of many conditions vary between different populations [35, 36, 37, 38]. Crucially, the features included in a dataset may be more or less useful for predicting different outcomes (e.g., being clinically diagnosed with a condition or not). The analysis of a dataset under this perspective is known as *data complexity* analysis, and it encompasses multiple aspects. A significant body of research has been dedicated to the formulation of various metrics that encapsulate the multifaceted aspects of dataset complexity [39]. Beyond disparities in model accuracy and data complexity, recent work highlights the importance of variance in model predictions. This variance is related to the extent to which model predictions can “flip” under minor changes in the training data, and it becomes an aspect of algorithmic fairness when high-variance predictions are concentrated in a demographic subgroup. This is called *systematic arbitrariness* [40]

This paper describes a multifaceted analysis of training datasets pertinent to chronic diseases aimed at uncovering potential discrepancies that could lead to biases in the resulting ML models. Our research substantiates the premise that demographic parity within datasets does not inherently ensure uniformity in algorithmic performance. That is to say, even datasets that are ostensibly equitable in terms of demographic attributes may still yield models with performance discrepancies. Initiating our analysis with a common ML performance metric, the Area Under the Receiving Operating Characteristic curve (AUROC or AUC), we measure the predictive efficacy of the models. Subsequently, our examination extends to more profound dataset attributes impacting model behavior, particularly data complexity and systematic arbitrariness. Our methodology provides a comprehensive approach for the assessment of training data from the perspective of algorithmic fairness. To the best of our knowledge, this investigation is the first to test systematic model arbitrariness in the healthcare domain.

## Methods

### Datasets

We use a list of datasets identified and reported in a survey of publicly accessible datasets related to chronic diseases [41]. Within this selection, two datasets pertain to diabetes ( $D_1, D_2$ ), while five are related to cardiac conditions ( $D_3 \dots D_7$ ). Dataset sizes vary widely (see Table 1), and for the purpose of this study, we segmented two of the large datasets into smaller subsets ( $D_{2a}, D_{2b}, D_{7a}, D_{7b}$ ) by randomly selecting two samples, each sized 100 times larger than the number of attributes. For the purpose of analysis, sex and age variables are binarized. In the case of age, the individuals within the lowest two quintiles are categorized as “young”, and those within the highest two quintiles are categorized as “old”, with the median quintile remaining unassigned. Dataset  $D_1$  does not include sex. Dataset  $D_{7\{a,b\}}$  was made available in 2020, but the specific year of data collection is not explicitly documented.

Table 1: Characteristics of the datasets used in this research.

| Dataset ID | Therapeutic Area | N      | Year    | Sex Ratio Female:Male | Younger Group Age Range | Elder Group Age Range |
|------------|------------------|--------|---------|-----------------------|-------------------------|-----------------------|
| $D_1$      | Diabetes         | 768    | 1988    | -                     | [21, 23]                | [33, 81]              |
| $D_{2a}$   | Diabetes         | 4,400  | 2014    | 1.17                  | [5, 65]                 | [75, 95]              |
| $D_{2b}$   | Diabetes         | 4,400  | 2014    | 1.09                  | [5, 65]                 | [75, 95]              |
| $D_3$      | Heart Dis.       | 920    | 1989    | 0.29                  | [28, 52]                | [57, 77]              |
| $D_4$      | Heart Dis.       | 452    | 1997    | 1.27                  | [0, 43]                 | [51, 83]              |
| $D_5$      | Heart Dis.       | 4,240  | 2010    | 1.33                  | [32, 46]                | [52, 70]              |
| $D_6$      | Heart Dis.       | 10,000 | 2020    | 0.79                  | [4, 57]                 | [66, 98]              |
| $D_{7a}$   | Heart Dis.       | 1,300  | ca.2020 | 1.86                  | [30, 52]                | [56, 65]              |
| $D_{7b}$   | Heart Dis.       | 1,300  | ca.2020 | 1.93                  | [30, 52]                | [56, 65]              |

## Model Performance

For the evaluation of model performance, we used three gradient boosting algorithms (XGBoost [42], LGBost [43], HGBost [44]) that support missing values. We considered two sets of attributes: including the protected attributes (“aware model”), and excluding them (“unaware model”). The performance metrics were similar across both models, which means that the datasets contain proxy variables for the protected attributes. Model training was done using a 3-fold cross validation schema, which involves partitioning the dataset into three subsets and cyclically using two-thirds for training and one-third for testing. This evaluation was further complemented by repeated bootstrapping, wherein each iteration involved a novel partitioning of the dataset. Hence, each reported Area Under ROC Curve value (AUROC, or simply ROC) is the average of 66 models: 3 algorithms times 22 runs (19 random runs plus 3 cross-validation runs).

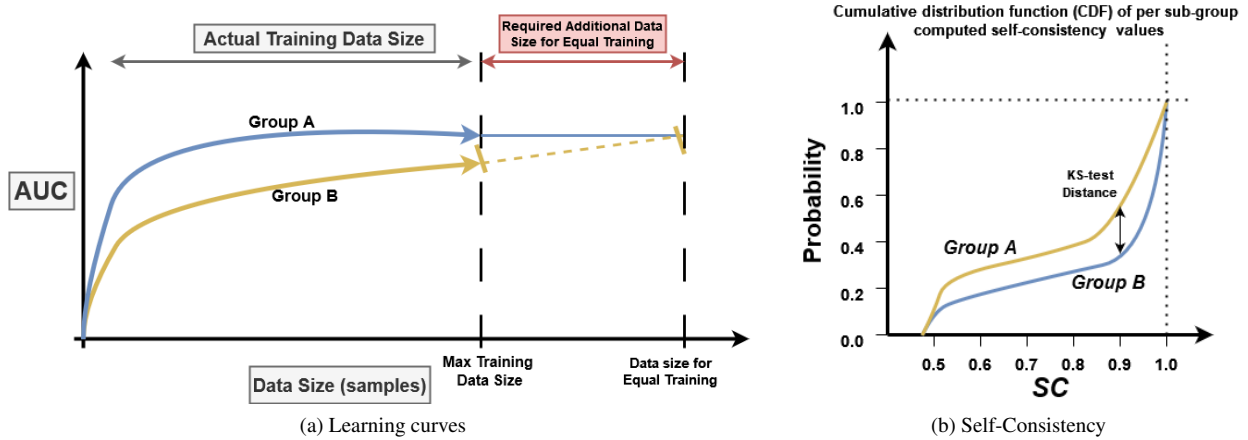


Figure 1: Depiction of our methods regarding learning curves and self-consistency.

Learning curves were obtained through an analogous process. We extrapolated learning curves to deduce an estimate of the number of additional data points that would enable the group with lower performance to attain the benchmark set by the group with higher performance. Figure 1a illustrates our method. Let  $f(n)$  be the superior learning curve, and  $g(n)$  be the inferior learning curve, with  $h(n)$  being an extrapolation of the learning curve. Conservatively, if we assume the upper curve reaches a saturation point  $f(N_p) = AUC_p$  (which is not always the case, hence the conservative estimate), we attempt the following minimization:

$$\begin{aligned} \min \quad & N_{add} \\ \text{s.t.} \quad & h(N_p + N_{add}) = f(N_p) \end{aligned}$$

i.e., we calculate the minimal number of additional data points that would be required from the group with the lower AUC to match the AUC of the group with higher performance. We consider three different functions  $h_1(\cdot)$ ,  $h_2(\cdot)$ ,  $h_3(\cdot)$ , each linearly constructed based on different segments of the concluding portion of the learning curve, choosing the one that necessitates the smallest increase in data points (i.e., the most conservative scenario, to avoid exaggerating the discrepancy). The greater the value of  $N_{add}$ , the larger the performance disparity.

## Data Complexity Metrics

Data complexity analysis is a systematic effort to understand discrepancies in classification accuracy by relating them to intrinsic characteristics of a dataset. This is a large research topic, and the interested reader can consult any of various surveys about it [45, 46, 47]. We used a fairly standard categorization of data complexity metrics [39], and picked one popular complexity metric within each category, as shown in Table 2.

Table 2: Data complexity: metrics categories and representative used; see [39] for details.

| Family name     | Object of analysis           | Metric used                                     |
|-----------------|------------------------------|---|
| Feature-based   | Feature informativeness      | Max. Fisher’s discriminant ratio                |
| Linearity       | Linear separability          | Sum of the error distance by linear programming |
| Neighborhood    | Local class distribution     | Error rate of nearest-neighbors classifier      |
| Dimensionality  | Data dimensionality/sparsity | Average number of features per point            |
| Class imbalance | Ratio between class examples | Imbalance ratio                                 |

### *Systematic Arbitrariness*

A family of models (e.g., various models built using the same learning scheme but different portions of the training data) may exhibit arbitrariness. Model arbitrariness corresponds to discrepancies in the predicted label for some elements across models of the same family, and it tends to be systematic, i.e., concentrated on specific items.

A recent study introduces a metric of Self-Consistency (SC) [40], which is computed at the level of an item as the probability that two models of the same family agree on the label for an item. For instance, an item with self-consistency of 1.0 is an item for which any model of a family yields the same predicted label. In binary classification, the minimum self-consistency is 0.5, indicating that half of the models yield one predicted label, and half of the models yield the opposite label. Note that self-consistency is independent of the “true” label of an element.

To compare self-consistency scores between groups, as recommended in [40] we use the Cumulative Distribution Function (CDF) of self-consistency. Often, one curve is above another, similarly to what we see in Figure 1b. We measure the disparity by performing a statistical test of the difference between the two curves.

## **Results**

### *Variations in model performance*

Using datasets collected in prior research, which include patient demographic details such as sex and age [41], our approach leverages three distinct gradient boosting algorithms to infer ML models from the training data. The validation methodology used herein incorporates cross-validation complemented by iterative bootstrapping, thereby generating a multitude of models each informed by different subsets of the training data. By examining the Area Under the ROC Curve (AUC), we determine the model’s proficiency in distinguishing between patient cohorts with different clinical outcomes. Table 3 presents a disaggregated view of AUC discrepancies across sex and age demographics. These results account for models that incorporate age and sex as predictive attributes (“aware modeling”). Similar results are observed when these variables are omitted (“unaware modeling”).

In our analysis, we observe disparities between sexes within several models, and across age groups in all but one model. Regarding sex-based disparities, approximately 10% of validation results reveal a higher AUC for males compared to females, whereas a mere 1% of results show higher female AUC relative to male AUC. Regarding age-related variances, these disparities exceed the sex-related ones with 32% of validation results demonstrating that the AUC for younger patients exceeds that of older patients, and conversely, in 5% of the cases, the AUC is greater for older patients compared to the younger patients.

### *Learning curves and the expected impact of additional data*

Learning curves are a standard tool for monitoring changes in model performance with the incremental addition of training data points. These curves graphically depict a performance metric, such as AUC, against the volume of training data utilized to build the model. Through this visual representation, one can appreciate trends like the

Table 3: Model performance (AUC), with average and variance computed over 66 models for each cell in the table. Differences along sex/age are expressed using p-values:  $< 0.01$  (\*),  $< 0.001$  (\*\*),  $< 0.0001$  (\*\*\*). The highest AUC is in boldface when the difference is significant at  $p < 0.01$ .

| Dataset  | AUC                               |                                   |                |                                   |                                   |                |
|----------|-----------------------------------|-----------------------------------|----------------|-----------------------------------|-----------------------------------|----------------|
|          | Female                            | Male                              | p-value        | Old                               | Young                             | p-value        |
| $D_1$    | –                                 | –                                 | –              | $0.65 \pm 0.04$                   | <b><math>0.69 \pm 0.05</math></b> | $< 0.0001$ *** |
| $D_{2a}$ | $0.59 \pm 0.02$                   | <b><math>0.61 \pm 0.02</math></b> | $< 0.0001$ *** | $0.58 \pm 0.02$                   | <b><math>0.62 \pm 0.02</math></b> | $< 0.0001$ *** |
| $D_{2b}$ | $0.59 \pm 0.02$                   | $0.59 \pm 0.02$                   | 0.67           | $0.58 \pm 0.02$                   | <b><math>0.60 \pm 0.02</math></b> | $< 0.0001$ *** |
| $D_3$    | $0.66 \pm 0.07$                   | <b><math>0.71 \pm 0.06</math></b> | $< 0.0001$ *** | $0.63 \pm 0.06$                   | <b><math>0.73 \pm 0.06</math></b> | $< 0.0001$ *** |
| $D_4$    | $0.77 \pm 0.04$                   | $0.78 \pm 0.06$                   | 0.80           | $0.79 \pm 0.05$                   | $0.79 \pm 0.06$                   | 0.95           |
| $D_5$    | $0.53 \pm 0.02$                   | <b><math>0.56 \pm 0.02</math></b> | $< 0.0001$ *** | <b><math>0.54 \pm 0.02</math></b> | $0.50 \pm 0.01$                   | $< 0.0001$ *** |
| $D_6$    | <b><math>0.92 \pm 0.01</math></b> | $0.87 \pm 0.02$                   | $< 0.0001$ *** | $0.87 \pm 0.02$                   | <b><math>0.91 \pm 0.01</math></b> | $< 0.0001$ *** |
| $D_{7a}$ | <b><math>0.68 \pm 0.03</math></b> | $0.67 \pm 0.03$                   | $< 0.01$ *     | $0.62 \pm 0.03$                   | <b><math>0.69 \pm 0.03</math></b> | $< 0.0001$ *** |
| $D_{7b}$ | <b><math>0.71 \pm 0.02</math></b> | $0.67 \pm 0.03$                   | $< 0.0001$ *** | $0.66 \pm 0.02$                   | <b><math>0.70 \pm 0.03</math></b> | $< 0.0001$ *** |

plateauing of performance gains and extrapolate the requisite quantity of additional data points necessary to attain a predefined AUC level in scenarios where the learning curve does not plateau.

Consistently aligning with our previous results regarding the AUC obtained from the comprehensive training datasets (excluding the fraction set aside for testing during cross-validation), it is often observed that the learning curve for one demographic group is above the learning curve for the other. This phenomenon suggests that even with balanced training sets, the resultant AUC may favor one group over another. Such a trend provides indirect evidence of differences in the predictability of outcomes between two groups when identical training data volumes and features are employed.

We have estimated the number of additional data points from the group exhibiting lower performance that would be required to match the AUC of the group with superior performance, based on the extrapolation of the learning curves to the point of anticipated AUC parity, given the current trajectory. Table 4 summarizes our results. The symbol Infinity ( $\infty$ ) means that the learning curve for the group for which the model has lower performance appears to plateau at an AUC threshold, signifying no further enhancement with additional data points.

Our analysis reveals that there are often imbalances that, to be corrected, would require a substantive amount of additional training instances. In half of the datasets, equating the AUC for sex would require the addition of 2% to 57% additional data for females, while for the remainder, an increment of 3% and 48% would be required for males. Regarding age, achieving parity would involve acquiring 2% and 46% more data for the younger patients in two datasets, and a substantial 5% and 192% increase for the older patients in the other datasets.

Table 4: Estimation, obtained by extrapolating learning curves, of the additional data points ( $N_{add}$ ) needed to achieve AUC parity.

| Dataset  | Group  | $N_{add}/N$ | Group | $N_{add}/N$ |
|----------|--------|-------------|-------|-------------|
| $D_1$    | –      | –           | Old   | 192%        |
| $D_{2a}$ | Female | 13%         | Old   | 112%        |
| $D_{2b}$ | Female | 2%          | Old   | 129%        |
| $D_3$    | Female | 66%         | Old   | $\infty$    |
| $D_4$    | Male   | 3%          | Young | 2%          |
| $D_5$    | Female | 57%         | Young | 46%         |
| $D_6$    | Male   | 48%         | Old   | 8%          |
| $D_{7a}$ | Male   | 6%          | Old   | 5%          |
| $D_{7b}$ | Male   | 4%          | Old   | 33%         |

Focusing on datasets that would benefit from at least a 10% increase in data, in 3 out of the 8 datasets where sex data is present, additional training data are required for females. In a similar vein, for age-related imbalances, 4 out of the 9 datasets would need additional data for the older patients group to achieve AUC parity.

### Alignment of data complexity with some disparities

We considered sixteen complexity metrics grouped into five categories, each corresponding to a unique conceptual framework, and computed the disparity of each metric between protected subgroups regarding sex and age. For each data set, AUC disparity divided by complexity metric disparity creates a ratio reflecting the consistency between model performance and data complexity as follows (where CM: Complexity Metric and A,B: Sub-groups A and B, i.e. Female-Male and Old-Young):

$$\frac{\overline{\text{AUC}}_A - \overline{\text{AUC}}_B}{\text{CM}_B - \text{CM}_A} = \begin{cases} 1 \text{ (Consistency)} & \text{if } x > 0 \\ -1 \text{ (Inconsistency)} & \text{if } x \leq 0 \end{cases}$$

Figure 2 presents the results in a heatmap visualization highlighting with light color the cases in which higher complexity and lower AUC values are observed for a specific sub-group in comparison with the other, while dark color indicates inconsistent AUC and complexity patterns.

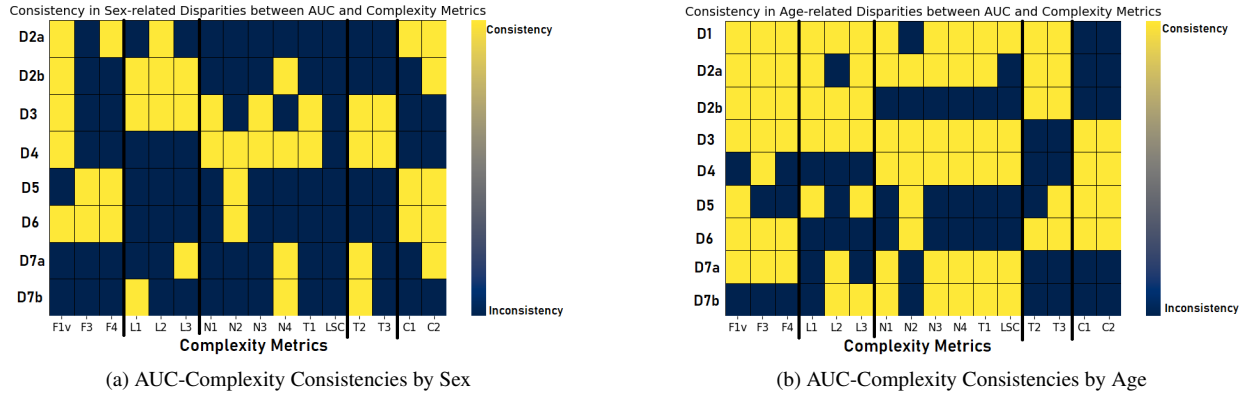


Figure 2: Consistency of sex- and age-related disparities between AUC and complexity metrics. Light colors indicate the cases in which auc and data complexity disparities are consistent, while dark colors inconsistent.

No obvious patterns can be observed across data in the results. Indeed, there are some situations of complementarity in which complexity metrics that are well aligned with AUC in some datasets are not aligned with AUC for another dataset and vice versa. Nevertheless, more complex data could potentially be linked to lower model performance, as homogeneous behavior is observed for some categories of metrics (especially in feature, dimensionality and class imbalance) within datasets regarding age. In addition, several sets of databases (e.g. those related to diabetes  $D_1, D_{2a}, D_{2b}$ ) show consistent disparities between performance and i) feature-based and ii) dimensionality complexity metrics. However, these experiments suggest that complexity metrics cannot be relied upon as a predictor of AUC disparities in specific clinical conditions.

### Systematic arbitrariness and model stability

In analyzing a family of models, each trained on distinct yet equivalently-sized partitions of the training data, we define an individual's self-consistency as the probability that two models within this family will yield the same label [40]. In our case, the minimum self-consistency is attained by subjects for which half of the models predict that they will be diagnosed with a condition, while the other half predict that they will not. Evidently, this is a situation we would like to avoid as much as possible. Hence, all other things equal, a model with higher self-consistency for most items is preferable.

Systematic arbitrariness is observed when items with low self-consistency are concentrated within a particular group, and can be measured by comparing CDFs. To quantify disparities, we use the Kolmogorov-Smirnov (KS) statistical test (Figure 1b), with the results shown in Table 5.

Table 5: Sub-group arbitrariness: area under the CDF of self-consistency results for each sub-group and distance between them measured by Kolmogorov-Smirnov (KS) statistical test. The results of the significance test indicate  $p$ -value  $< 0.01$  (\*),  $< 0.001$  (\*\*), and when significant the larger arbitrariness appears in boldface.

| Dataset  | Overall | Sex    |       |         | Age          |         |         | KS-test |
|----------|---------|--------|-------|---------|--------------|---------|---------|---------|
|          |         | Female | Male  | KS-test | Elder        | Younger | KS-test |         |
| $D_1$    | 0.243   | -      | -     | -       | <b>0.247</b> | 0.172   | 0.205   | *       |
| $D_{2a}$ | 0.323   | 0.308  | 0.312 | 0.019   | 0.295        | 0.294   | 0.045   |         |
| $D_{2b}$ | 0.332   | 0.313  | 0.319 | 0.020   | 0.302        | 0.299   | 0.018   |         |
| $D_3$    | 0.258   | 0.223  | 0.238 | 0.071   | 0.222        | 0.210   | 0.043   |         |
| $D_4$    | 0.313   | 0.233  | 0.247 | 0.037   | 0.279        | 0.263   | 0.057   |         |
| $D_5$    | 0.119   | 0.105  | 0.119 | 0.060   | <b>0.169</b> | 0.069   | 0.299   | **      |
| $D_6$    | 0.084   | 0.083  | 0.078 | 0.031   | 0.078        | 0.085   | 0.024   |         |
| $D_{7a}$ | 0.266   | 0.244  | 0.233 | 0.053   | <b>0.263</b> | 0.233   | 0.082   | *       |
| $D_{7b}$ | 0.269   | 0.247  | 0.250 | 0.049   | <b>0.262</b> | 0.228   | 0.088   | *       |

Results show that ML predictions have no significant self-consistency disparities between male and female subjects. However, older individuals exhibit significantly more arbitrary predictions in 4 out of 9 datasets. These results are for a model including protected characteristics (“unaware model”). Results for the aware model are similar, although with lower self-consistency values in general and smaller differences by sex and age.

## Discussion

Our findings uncover sex- and age-related disparities in model performance as evidenced by the AUC of the models. It is pertinent to recall that the representation of each group within the datasets is equal. This observation underscores that mere demographic parity in training datasets does not mean model equity. The analysis of learning curves provides insights into the potential benefits of data augmentation. Sex-related disparities are observed to occasionally favor males over females, with a marginal predominance for male patients as indicated by both AUC and the requisite additional data to attain performance parity. Regarding age differences, the findings are more pronounced, with models generally predicting better for younger patients across most datasets (higher AUC), and requiring a large volume of additional training data to potentially achieve performance parity.

Furthermore, upon examining disparities in data complexity and systematic arbitrariness, we observe that predictions for older patients tend to be less consistent than those for their younger counterparts in several datasets. These disparities, to some extent, correlate with the model performance (AUC) and data complexity findings, suggesting a linkage between increased data arbitrariness for older patients and heightened complexity, leading to lower model performance. These correlations suggest but do not determine model disparities, as there are exceptions within our observations, where greater arbitrariness is sometimes associated with comparable or superior AUC values. This highlights the necessity of a multifaceted metric consideration encompassing performance, complexity, and stability, rather than relying exclusively on performance metrics.

Within the healthcare domain, the legal and ethical dimensions of decision-making are of paramount importance [48]. The findings of this study highlight some characteristics of model performance that are not typically reported, but that hold considerable potential to influence clinical practice. Specifically, systematic arbitrariness in model outputs could undermine clinician confidence in ML and diminish the acceptability of such models. We propose datasets are tested for systematic arbitrariness before being used in clinical settings. In nearly half of the datasets we studied, older patients with chronic diseases face the risk of health inequities [49, 50] due to data that is suboptimal for modeling their health outcomes as compared to younger patients.

Hospital data, such as the one used in this study, may be indicative solely of the population with healthcare system access, thus potentially engendering bias against certain subpopulations [51, 52, 53, 54]. Future efforts should aim to extend these analyses to include additional databases. To address situations where systematic arbitrariness is detected, we must consider both technical and human factors [55]. This includes designing systems that minimize potential technology-induced disparities, taking into account the data and algorithmic literacy of the users of these systems, i.e., clinicians. Arbitrariness is not a new concept in the health domain, as evidenced by the existence of cost-effective pharmacological treatments that exhibit suboptimal efficacy in particular patient subgroups. Systems are in place

to educate and safeguard against potential patient harm, including rigorous and multiphase pharmaceutical clinical studies and pharmacovigilance protocols. Data quality audits should scrutinize performance differentials impacting specific subgroups, whose data characteristics may differ from other populations with the same condition. Future research could explore the creation of monitoring processes for ML models in healthcare, analogous to those applied to pharmacological drugs.

## Conclusion

In this study, we identified significant age-related and mild sex-related disparities in the performance of ML models for chronic disease prediction. Older patients, in particular, experienced inconsistent and arbitrary predictions across several datasets due to increased data complexity and lower model performance, while sex-based differences slightly favored male predictions. These findings demonstrate that representativeness in training data alone is insufficient for ensuring equitable outcomes. Therefore, addressing model arbitrariness, especially for older individuals, is essential before deploying ML models in clinical settings to ensure fairness and reliability.

**Acknowledgments** This work is funded by Spanish Ministry of Science and Innovation (REF:DIN2021-011865). In addition, it has been partially supported by the Department of Research and Universities of the Government of Catalonia (SGR 00930), EU-funded project "SoBigData++" (grant agreement 871042), and MCIN/AEI/10.13039/501100011033 under the Maria de Maeztu Units of Excellence Programme (CEX2021-001195-M). None of the funders played any role in study design, data collection, analysis and interpretation of data, or the writing of this manuscript.

**Data Availability** A version of the datasets used during the current study that has been pre-processed for reproducibility is available in the following repository, including pointers to the original datasets:  
[https://anonymous.4open.science/r/health\\_disparities\\_Chronic\\_Diseases-32D4/README.md](https://anonymous.4open.science/r/health_disparities_Chronic_Diseases-32D4/README.md)

## References

1. Fogel AL, Kvedar JC. Artificial intelligence powers digital medicine. *NPJ digital medicine*. 2018;1(1):5.
2. Jusman Y, Firdiantika IM, Dharmawan DA, Purwanto K. Performance of multi layer perceptron and deep neural networks in skin cancer classification. In: 2021 IEEE 3rd global conference on life sciences and technologies (LifeTech). IEEE; 2021. p. 534-8.
3. Lo B, Malina D, Pittman G, Morrissey S. *Fundamentals of medical ethics—A new perspective series*. Mass Medical Soc; 2023.
4. Hamam H. Guest editorial achieving health equity through AI for diagnosis and treatment and patient monitoring. *IEEE Journal of Biomedical and Health Informatics*. 2024;28(2):702-6.
5. Liu M, Ning Y, Teixayavong S, Mertens M, Xu J, Ting DSW, et al. A translational perspective towards clinical AI fairness. *NPJ Digital Medicine*. 2023;6(1):172.
6. Fletcher RR, Nakeshimana A, Olubeko O. Addressing fairness, bias, and appropriate use of artificial intelligence and machine learning in global health. *Frontiers Media SA*; 2021.
7. Buolamwini J, Geburu T. Gender shades: Intersectional accuracy disparities in commercial gender classification. In: *Conference on fairness, accountability and transparency*. PMLR; 2018. p. 77-91.
8. Hajian S, Bonchi F, Castillo C. Algorithmic bias: From discrimination discovery to fairness-aware data mining. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*; 2016. p. 2125-6.
9. Raza S. Connecting fairness in machine learning with public health equity. In: 2023 IEEE 11th International Conference on Healthcare Informatics (ICHI). IEEE; 2023. p. 704-8.
10. Lippert-Rasmussen K. *Born free and equal?: A philosophical inquiry into the nature of discrimination*. Oxford University Press; 2013.
11. Feuerriegel S, Dolata M, Schwabe G. Fair AI: Challenges and opportunities. *Business & information systems engineering*. 2020;62:379-84.
12. Dastin J. Amazon scraps secret AI recruiting tool that showed bias against women. In: *Ethics of data and analytics*. Auerbach Publications; 2022. p. 296-9.
13. Caliskan A, Bryson JJ, Narayanan A. Semantics derived automatically from language corpora contain human-like biases. *Science*. 2017;356(6334):183-6.



14. Hardesty L. Study finds gender and skin-type bias in commercial artificial-intelligence systems. Retrieved April. 2018;3:2019.
15. Tsai TC, Arik S, Jacobson BH, Yoon J, Yoder N, Sava D, et al. Algorithmic fairness in pandemic forecasting: lessons from COVID-19. *NPJ Digital Medicine*. 2022;5(1):59.
16. Halamka J, Bydon M, Cerrato P, Bhagra A. Addressing racial disparities in surgical care with machine learning. *NPJ Digital Medicine*. 2022;5(1):152.
17. Drukker K, Chen W, Gichoya J, Grusauskas N, Kalpathy-Cramer J, Koyejo S, et al. Toward fairness in artificial intelligence for medical image analysis: identification and mitigation of potential biases in the roadmap from data collection to model deployment. *Journal of Medical Imaging*. 2023;10(6):061104-4.
18. Larrazabal AJ, Nieto N, Peterson V, Milone DH, Ferrante E. Gender imbalance in medical imaging datasets produces biased classifiers for computer-aided diagnosis. *Proceedings of the National Academy of Sciences*. 2020;117(23):12592-4.
19. Seyyed-Kalantari L, Liu G, McDermott M, Chen IY, Ghassemi M. CheXclusion: Fairness gaps in deep chest X-ray classifiers. In: *BIOCOMPUTING 2021: proceedings of the Pacific symposium*. World Scientific; 2020. p. 232-43.
20. Chen IY, Szolovits P, Ghassemi M. Can AI help reduce disparities in general medical and mental health care? *AMA journal of ethics*. 2019;21(2):167-79.
21. Burlina P, Joshi N, Paul W, Pacheco KD, Bressler NM. Addressing artificial intelligence bias in retinal diagnostics. *Translational Vision Science & Technology*. 2021;10(2):13-3.
22. Abubakar A, Ugail H, Bukar AM. Assessment of human skin burns: a deep transfer learning approach. *Journal of Medical and Biological Engineering*. 2020;40:321-33.
23. Thompson HM, Sharma B, Bhalla S, Boley R, McCluskey C, Dligach D, et al. Bias and fairness assessment of a natural language processing opioid misuse classifier: detection and mitigation of electronic health record data disadvantages across racial subgroups. *Journal of the American Medical Informatics Association*. 2021;28(11):2393-403.
24. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*. 2019;366(6464):447-53.
25. Park Y, Hu J, Singh M, Sylla I, Dankwa-Mullan I, Koski E, et al. Comparison of methods to reduce bias from clinical prediction models of postpartum depression. *JAMA network open*. 2021;4(4):e213909-9.
26. Jobin A, Ienca M, Vayena E. The global landscape of AI ethics guidelines. *Nature machine intelligence*. 2019;1(9):389-99.
27. Marabelli M, Newell S, Handunge V. The lifecycle of algorithmic decision-making systems: Organizational choices and ethical challenges. *The Journal of Strategic Information Systems*. 2021;30(3):101683.
28. Leavy S. Gender bias in artificial intelligence: The need for diversity and gender theory in machine learning. In: *Proceedings of the 1st international workshop on gender equality in software engineering*; 2018. p. 14-6.
29. Anderson JW, Shaikh N, Visweswaran S. Measuring and reducing racial bias in a pediatric urinary tract infection model. *AMIA Summits on Translational Science Proceedings*. 2024;2024:488.
30. Suresh H, Guttig J. A framework for understanding sources of harm throughout the machine learning life cycle. In: *Equity and access in algorithms, mechanisms, and optimization*; 2021. p. 1-9.
31. Ahmad MA, Patel A, Eckert C, Kumar V, Teredesai A. Fairness in machine learning for healthcare. In: *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*; 2020. p. 3529-30.
32. Corbett-Davies S, Pierson E, Feller A, Goel S, Huq A. Algorithmic decision making and the cost of fairness. In: *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*; 2017. p. 797-806.
33. Raji ID, Buolamwini J. Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial ai products. In: *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*; 2019. p. 429-35.
34. Yang Y, Gupta A, Feng J, Singhal P, Yadav V, Wu Y, et al. Enhancing fairness in face detection in computer vision systems by demographic bias mitigation. In: *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*; 2022. p. 813-22.
35. Boehme AK, Siegler JE, Mullen MT, Albright KC, Lyerly MJ, Monlezun DJ, et al. Racial and gender differences in stroke severity, outcomes, and treatment in patients with acute ischemic stroke. *Journal of Stroke and Cerebrovascular Diseases*. 2014;23(4):e255-61.

36. Cirillo D, Catuara-Solarz S, Morey C, Guney E, Subirats L, Mellino S, et al. Sex and gender differences and biases in artificial intelligence for biomedicine and healthcare. *NPJ digital medicine*. 2020;3(1):81.
37. Celeste C, Ming D, Broce J, Ojo DP, Drobina E, Louis-Jacques AF, et al. Ethnic disparity in diagnosing asymptomatic bacterial vaginosis using machine learning. *NPJ Digital Medicine*. 2023;6(1):211.
38. Hou B, Mondragón A, Tarzanagh DA, Zhou Z, Saykin AJ, Moore JH, et al. PFERM: A fair empirical risk minimization approach with prior knowledge. *AMIA Summits on Translational Science Proceedings*. 2024;2024:211.
39. Lorena AC, Garcia LP, Lehmann J, Souto MC, Ho TK. How complex is your classification problem? a survey on measuring classification complexity. *ACM Computing Surveys (CSUR)*. 2019;52(5):1-34.
40. Cooper AF, Lee K, Choksi MZ, Barocas S, De Sa C, Grimmelmann J, et al. Arbitrariness and social prediction: The confounding role of variance in fair classification. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 38; 2024. p. 22004-12.
41. Bilonis I, Fernandez-Luque L, Castillo C. A survey on public data sets related to chronic diseases. In: *2023 IEEE 36th International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE; 2023. p. 917-20.
42. Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*; 2016. p. 785-94.
43. Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, et al. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*. 2017;30.
44. Guryanov A. Histogram-based algorithm for building gradient boosting ensembles of piecewise linear decision trees. In: *Analysis of Images, Social Networks and Texts: 8th International Conference, AIST 2019, Kazan, Russia, July 17–19, 2019, Revised Selected Papers 8*. Springer; 2019. p. 39-50.
45. Santos MS, Abreu PH, Japkowicz N, Fernández A, Santos J. A unifying view of class overlap and imbalance: Key concepts, multi-view panorama, and open avenues for research. *Information Fusion*. 2023;89:228-53.
46. Rivolli A, Garcia LP, Soares C, Vanschoren J, de Carvalho AC. Meta-features for meta-learning. *Knowledge-Based Systems*. 2022;240:108101.
47. Gupta N, Mujumdar S, Patel H, Masuda S, Panwar N, Bandyopadhyay S, et al. Data quality for machine learning tasks. In: *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*; 2021. p. 4040-1.
48. Peters D, Vold K, Robinson D, Calvo RA. Responsible AI—two frameworks for ethical design practice. *IEEE Transactions on Technology and Society*. 2020;1(1):34-47.
49. Lorenc T, Petticrew M, Welch V, Tugwell P. What types of interventions generate inequalities? Evidence from systematic reviews. *J Epidemiol Community Health*. 2012.
50. Klang E, Soffer S, Nadkarni G, Glicksberg B, Freeman R, Horowitz C, et al. Sex differences in age and comorbidities for COVID-19 mortality in urban New York City. *SN comprehensive clinical medicine*. 2020;2:1319-22.
51. Chen D, Liu S, Kingsbury P, Sohn S, Storlie CB, Habermann EB, et al. Deep learning and alternative learning strategies for retrospective real-world clinical data. *NPJ digital medicine*. 2019;2(1):43.
52. Hyun KK, Redfern J, Patel A, Peiris D, Brieger D, Sullivan D, et al. Gender inequalities in cardiovascular risk factor assessment and management in primary healthcare. *Heart*. 2017;103(7):492-8.
53. Gauci S, Cartledge S, Redfern J, Gallagher R, Huxley R, Lee CMY, et al. Biology, bias, or both? The contribution of sex and gender to the disparity in cardiovascular outcomes between women and men. *Current Atherosclerosis Reports*. 2022;24(9):701-8.
54. Veinot TC, Mitchell H, Ancker JS. Good intentions are not enough: how informatics interventions can worsen inequality. *Journal of the American Medical Informatics Association*. 2018;25(8):1080-8.
55. Ferryman K, Mackintosh M, Ghassemi M. Considering biased data as informative artifacts in AI-assisted health care. *New England Journal of Medicine*. 2023;389(9):833-8.