# Uncovering Bias in Personal Informatics

SOFIA YFANTIDOU, PAVLOS SERMPEZIS, and ATHENA VAKALI, Aristotle University of Thessaloniki, Greece

RICARDO BAEZA-YATES, Institute for Experiential AI, Northeastern University, United States

Personal informatics (PI) systems, powered by smartphones and wearables, enable people to lead healthier lifestyles by providing meaningful and actionable insights that break down barriers between users and their health information. Today, such systems are used by billions of users for monitoring not only physical activity and sleep but also vital signs and women's and heart health, among others. Despite their widespread usage, the processing of sensitive PI data may suffer from biases, which may entail practical and ethical implications. In this work, we present the first comprehensive empirical and analytical study of bias in PI systems, including biases in raw data and in the entire machine learning life cycle. We use the most detailed framework to date for exploring the different sources of bias and find that biases exist both in the data generation and the model learning and implementation streams. According to our results, the most affected minority groups are users with health issues, such as diabetes, joint issues, and hypertension, and female users, whose data biases are propagated or even amplified by learning models, while intersectional biases can also be observed.

CCS Concepts: • **Human-centered computing** → **Ubiquitous and mobile computing**; • **Applied computing** → *Consumer health*; • **Computing methodologies** → *Artificial intelligence*; • **Social and professional topics** → *Codes of ethics*.

Additional Key Words and Phrases: machine learning, bias, fairness, personal informatics, ubiquitous computing, sensing data, digital biomarkers

## 1 INTRODUCTION

Ubiquitous technologies, such as smartphones or wearables, are an integral part of our lives today, with the current number of smartphone users worldwide in 2022 rising to 6.6 billion, or 83.4% of the global population, from 3.6 billion (49.4%) in 2016 [48]. More impressively, the number of connected wearable devices worldwide has more than tripled in the past six years, rising to 1.1 billion in 2022 from 325 million in 2016 [92]. The proliferation of ubiquitous technologies has given rise to Personal Informatics (PI), namely a class of systems that "help people collect personally relevant information for the purpose of self-reflection and gaining self-knowledge" [67]. Such systems enable people to keep track of their productivity [61], finances [59], and learning [47]. Yet, tracking various aspects of physical and mental health is particularly prevalent [37].

Ubiquitous technologies for PI can continuously and unobtrusively measure and collect physiological and behavioral data, namely, "digital biomarkers", from users through integrated sensors. Digital biomarkers contain an uncanny amount of personal information. Even the coarser behavioral biomarkers acquired from consumer wearable devices (such as steps, burned calories, and covered distance), strongly correlate to a person's gender, height, and weight [60], while signals of finer granularity (such as accelerometer and heart rate measurements), can predict variables associated with an individual's physical health, fitness, and demographics [91]. Similarly,

Authors' addresses: Sofia Yfantidou, syfantid@csd.auth.gr; Pavlos Sermpezis, sermpezis@csd.auth.gr; Athena Vakali, avakali@csd.auth.gr, Aristotle University of Thessaloniki, Thessaloniki, Greece, 54124; Ricardo Baeza-Yates, Institute for Experiential AI, Northeastern University, San Jose, United States, CA 95113.

physiological signals accompanied by PI usage logs have been used to predict mood, stress and overall mental health [95]. At the same time, the advanced features for health tracking that are continuously integrated into consumer smartphones and wearables [13, 14, 70, 96] now enable advanced analytics, such as atrial fibrillation identification, fertility prediction, fall, and crash detection, sleep apnea warnings, and are paving the future of mHealth.

**Bias in PI.** However, the prevalent PI adoption embeds important challenges due to the questionable transparency and unexplored biases in the systems' algorithms. Bias in machine learning is a source of unfairness that can lead to harmful consequences, such as discrimination [72]. The Fairness, Accountability, and Transparency in Machine Learning (FAT/ML) community defines fairness as a principle that "ensures that algorithmic decisions do not create discriminatory or unjust impacts when comparing across different demographics (e.g., race, sex, etc.)" [16]. Fairness is an inexorably subjective and context-dependent notion and incorporates different metrics for different definitions, some of which are even mutually incompatible [44]. Contrary to the common belief that algorithmic decisions are objective and unbiased by definition, a machine learning model may actually be inherently unfair by learning, preserving, or even amplifying historical biases existent in the data [82]. Real-world cases of unfair machine learning models are, unfortunately, abundant. Examples can be drawn from criminal justice [7], hiring practices [30], ad targeting [90], facial recognition [84], healthcare [108] and language models [21].

Despite this growing interest in machine learning biases overall, a focused emphasis on the requirements of unbiased PI systems in mHealth settings is lacking [3]. PI systems are deployed in high-stakes health-related applications, while their input data modality makes them susceptible to propagating or even amplifying bias. Beyond algorithm performance, the existence of bias is a challenging problem in delivering equitable care. Thus, it is critical to explore biases within these systems to raise awareness regarding mitigating and regulatory actions required to avert potential negative consequences.

**PI Idiosyncracies.** Moreover, this need for exploring bias is further highlighted by the fact that the PI domain has significant differences -in terms of bias- from previously well-studied domains, such as facial or speech recognition:

- *The digital divide as a barrier of entry:* To contribute data to an image or voice dataset, users do not need any prerequisite knowledge or niche device. However, to contribute to a PI dataset, users face significant "entry barriers" in terms of digital capacity or device ownership, creating new-found *representation biases* in the domain's datasets, as verified by our analysis in Sections 3.1 and 3.2.
- *Emerging technologies accuracy:* Facial or speech recognition measurement devices, e.g., camera or voice recorder, are based on mature technologies. As a result, their accuracy remains relatively unchangeable across different devices. On the contrary, emerging PI devices' accuracy significantly varies across manufacturers and even across models, creating unexplored *measurement biases* and discrepancies between user segments (cf. Section 3.3).
- *Complex nature of data:* It may be easy to identify biases in terms of skin color and gender (facial recognition) or accent and gender (speech recognition). Yet, identifying biases in digital biomarkers (e.g., step or sleep data) may not be straightforward. Biases in PI data can remain hidden and be further propagated or even amplified in machine learning models (cf. Sections 4.1 and 4.2).

**Summary of contributions.** Motivated by these idiosyncrasies and the gap in the literature, in this paper, we present the first comprehensive study on bias in PI: We adopt the most complete framework to date for understanding sources of harm in the machine learning life cycle Suresh and Guttag [94], explore biases in the data generation and model and implementation streams, and validate them in a real-world, large-scale PI dataset. During this process, we examine the suitability of different fairness metrics for digital biomarkers, initiating a conversation within the community on how to approach biases within ubiquitous mHealth.

Specifically, our research questions (RQs) and the respective contributions of our study are as follows:

(1) ***What does bias mean for PI?*** To quantify bias in relation to the PI domain, we explore diverse fairness definitions and metrics and identify differences from other domains. We delineate each metric's strengths and shortcomings and select the most appropriate metrics for the domain (Section 2).

(2) ***Are PI data susceptible to biases?*** We examine the largest real-world PI dataset to date to assess whether ubiquitous digital biomarkers incorporate biases. Specifically, we perform the first detailed study on bias in the MyHeart Counts cardiovascular health study dataset [53], containing physical activity, fitness, sleep, and cardiovascular health data for 50K participants across the United States. Our results identify biases across all dimensions of the data generation stream, namely *historical*, *representation*, and *measurement* biases; these findings highlight that users should be cautious when using PI datasets, in general, and the MyHeart Counts data, in particular (Section 3).

(3) ***Do machine learning models inherit PI data biases?*** We investigate whether biases present in PI data are propagated when applying machine learning models to these data. Specifically, we evaluate long short-term memory (LSTM) sequence models as a baseline and personalized models for *aggregation*, *learning*, and *deployment* biases. In line with prior work [80], our findings indicate that data biases are propagated to deep learning models, especially for intersectional user groups. Surprisingly, they are significantly amplified in their personalized counterparts, raising questions regarding the shortcomings of personalization (Section 4).

(4) ***Can synthetic benchmarks hide the imperfect nature of PI?*** We explore whether "perfect" synthetic benchmark datasets can hide PI data and model "imperfections" and biases during evaluation. Specifically, we compare a random benchmark, representative of our data, with one designed to achieve demographic parity for evaluation biases. Our findings highlight the importance of the establishment of PI benchmarks that are representative of the intended target populations to avoid the deployment of models with unidentified biases (Section 4.3).

Finally, we partially apply our analysis on two different PI datasets to showcase the generalizability of our findings and share our code publicly [8] to encourage applicability to other datasets as well (Section 5). We then discuss limitations and recommendations for mitigating the effect of bias in PI and conclude our work (Section 7).
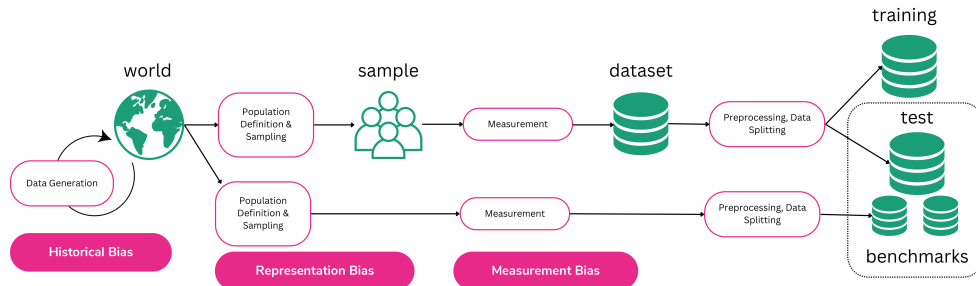
## 2  PERSONAL INFORMATICS BIASES: SETTING & CONFIGURATION

In this section, we discuss the framework upon which we base our study of bias in PI (Section 2.1) and describe our use case configuration, which acts as a starting point for our investigation (Section 2.2).
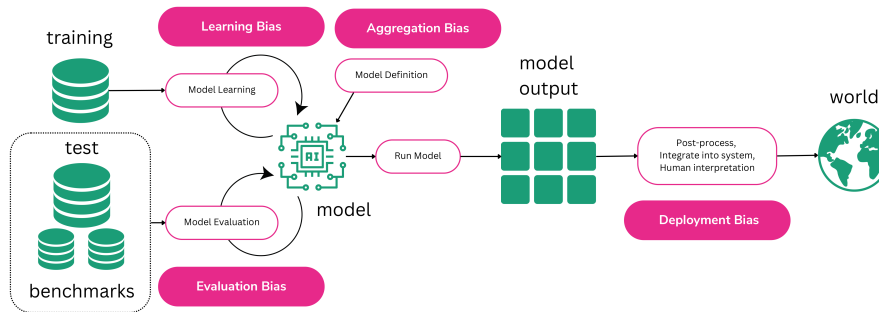
### 2.1  Sources of Bias in the Machine Learning Life Cycle Framework

We base our study of bias in PI for mHealth on Suresh and Guttag [94] framework for understanding sources of harm through the machine learning life cycle, the most comprehensive framework to date for capturing biases in any machine learning system. According to Suresh and Guttang, the machine learning life cycle consists of two major streams containing seven sources of bias-related harms, the *data generation* stream and the *model building and implementation* stream. The data generation stream can contain historical, representational, and measurement bias, while the model building and implementation stream can contain aggregation, learning, evaluation, and deployment bias. An overview of the sources of bias-related harm in the data generation and model building and implementation streams is shown in Figures 1a and 1b, respectively, while definitions are provided below [94]:

- *Historical biases* can occur even if the data are flawlessly measured and sampled by reflecting real-world, past, and present biases against one or more groups of people. For example, gender gaps in certain professions lead to natural language models associating gendered occupation words, such as nurse or programmer, with words representing women or men, respectively [21].

(a) Sources of harm in the data generation stream.



(b) Sources of harm in the model building and implementation stream.

Fig. 1. Sources of harm in the data (top) and model building and implementation (bottom) streams [94]. The training, test, and benchmark sets are common across figures.

- *Representation biases* can occur when sampling methods lead to underrepresenting general population segments. For example, in popular image datasets, the majority of the images originate from the United States or Europe, leading to performance degradation when classifying images coming from an underrepresented region [33].
- *Measurement biases* can occur when choosing, collecting, and calculating features and labels for the prediction problem. For example, in medical applications, oftentimes, diagnosis is used as a proxy for having a health condition; yet, certain gender and racial groups suffer higher rates of misdiagnosis, or underdiagnosis [54].
- *Aggregation biases* can occur when an "one-size-fits-all" treatment, e.g., model, is used for data in which underlying user groups should be treated separately. For example, in natural language processing, training models in generic data will fail to capture the different meanings, and off-line context behind street slang [43].
- *Learning biases* can occur when modeling choices amplify performance disparities across different user segments in the data. For example, optimizing a model for privacy can reduce the influence of data originating from underrepresented groups [18].
- *Evaluation biases* can occur when the benchmark population is not representative of the real user population. For example, dark-skinned women comprise only a small percentage of popular facial images benchmark, leading to worse performance of commercial facial analysis tools on intersectional accuracy [22].
- *Deployment biases* can occur when there exists a mismatch between the problem a model is designed to solve and how it is actually utilized. For example, risk assessment tools in criminal justice are not used in isolation but can be used in "off-label" ways, such as determining the length of a sentence [28].

In the following section, we introduce the use case through which we explore bias in PI for mHealth. We then show empirically and analytically how Suresh and Guttag's seven sources of bias translate in the PI domain.

## 2.2 Exploring Bias through the Largest Digital Biomarkers mHealth Dataset

To examine the existence of the seven sources of bias in the PI machine learning life cycle and provide clear answers to our research questions (RQs), introduced in Section 1, we need to define an indicative -but by no means restrictive- use case to enable our analysis. For this purpose, we utilize the MyHeart Counts dataset [53], the largest collection of digital biomarkers in the mHealth domain to date, enabling us to perform the most comprehensive analysis of bias across diverse user demographics, including gender, ethnicity, age, BMI, and health conditions. Nevertheless, our methodology and outputs can be generalized to any PI dataset other than the prominent use case of MyHeart Counts (see Section 5).

***Data Description***. Up till recently, general-purpose, population-scale PI datasets were unavailable, partly due to the high cost of data collection, as well as privacy concerns and data protection regulations. The most popular open datasets consisted of small to medium samples [99, 105] or were domain-constrained to Human-Activity Recognition (HAR) [6] and Sleep Classification (SC) [71]. However, this changed with the publication of data from the MyHeart Counts Cardiovascular Health Study, a collection of real-world physical activity, fitness, sleep, and cardiovascular health data from 50K participants in the United States. Participants completed various surveys and a 6-minute walk test and contributed PI data via an iPhone application built using Apple's ResearchKit framework [15]. They provided informed consent to make this data freely available for future research.

Table 1. The available protected attributes in the MyHeart Counts study data. For the purpose of the bias analysis, we convert the non-binary attributes to binary to ensure a sufficient sample size per group and compatibility with popular bias metrics.

| | *Original Protected Attribute Values* | *Binarized Protected Attribute Values* | |
|---|---|---|---|
| **Attribute** | **Original Groups** | **Majority Group** | **Minority Group** |
| Gender | Male, Female, N/A | Male | Female |
| Ethnicity | White, Asian, Black, Hispanic, American Indian, Pacific Islander, Other, N/A | White | Non-white |
| Age | Integer Number, N/A | <65 (lower risk of complications) | >=65 (higher risk of complications) |
| BMI | Real Number (height and weight), N/A | <18.5 or =>25 (non-healthy) | =>18.5 and <25 (healthy) |
| Heart Condition | Yes, No, N/A | No | Yes |
| Hypertension | Yes, No, N/A | No | Yes |
| Joint Problem | Yes, No, N/A | No | Yes |
| Diabetes | Yes, No, N/A | No | Yes |

Approximately 1 out of 10 participants ($N$ = 4920) shared their basic HealthKit data (step count, distance covered, burned calories, and flights climbed), while fewer users shared their sleep ($N$ = 626) and workout ($N$ = 881) data. We perform our analysis on the basic HealthKit data, which contains the most common data types among scientific datasets so that our findings are generalizable and our methodology is reproducible. Additionally, we combine these data with survey responses to attain the following user attributes: gender, ethnicity, age, BMI, and health conditions, such as heart condition, hypertension, joint problem, and diabetes.

***Data Preprocessing***. To ensure a sufficient sample size per user group and compatibility with popular bias metrics, we convert non-binary user attributes, such as ethnicity, age, and BMI, to binary, as seen in Table 1. This grouping creates two user groups per protected attribute, namely a majority group (also called "privileged" for the purpose of this analysis) and a minority group (also called "unprivileged" for the purpose of this analysis). Note that the usage of the term "privilege" in this work does not necessarily coincide with real-world "privilege". For example, users with non-healthy BMI are the majority user segment in our dataset, and hence, they are referred to as the "privileged" user group, whereas one could argue that the opposite applies in reality.

Table 2. An example of input data for the physical activity prediction use case. The step counts per hour for the past 48 hours are the features, and the total number of the next day's steps is the label. The user ID and timestamps are not used in the learning.

| | | Features | | | Label |
|---|---|---|---|---|---|
| user_id | timestamp | steps at t-48h | ... | steps at t-1h | next day's steps |
| 1 | 23-11-2022 | 1040 | ... | 300 | 8500 |

***Data Labeling***. As mentioned previously, the MyHeart Counts dataset is general-purpose, meaning that it does not introduce any new learning tasks or certain prediction labels. To this end, we select the *next-day physical activity prediction from historical data* use case [19, 101] for model training. In other words, based on the user's past activity, we try to predict how many steps they will perform the next day (see Table 2). Such a task may enable, for instance, the provision of personalized step goals by PI systems, which have proven to be more effective in inciting positive health behavior change compared to static, fixed goals [68]. The reasons behind this choice lie not only in the benefits of physical activity for physical and mental health [78] but also in the availability of basic digital behavioral biomarkers, such as steps. Contrary to raw sensor data, which are harder to collect at scale through consumer PI systems, basic digital behavioral biomarkers, are easy to collect and commonplace in the literature, enabling the reproducibility of our findings. At the same time, steps are the largest available sensed modality in the My Heart Counts dataset, allowing us to take advantage of a larger portion of the data for the purpose of this analysis. Finally, it is important to note that our findings can be generalized to other PI tasks, e.g., mood and stress prediction [95], or health monitoring [83].

## 2.3 Quantifying Bias through Fairness Metrics

In this section, we discuss fairness metrics that quantify bias in machine learning from the perspective of PI, providing an answer to RQ1, namely: *What does bias mean for PI?*

As mentioned earlier on, fairness is a social construct that defies simple definition [76]. Quantitative fields view fairness as a mathematical problem of "equal or equitable allocation, representation, or error rates, for a particular task or problem" [76]. There is a variety of fairness definitions and metrics (see Appendix A). However, not all of them are relevant to the PI domain. Contrary to popular bias quantification tasks, such as recidivism prediction or loan repayment prediction, in our use case -and related PI tasks- there is no clear positive outcome for the user. In other words, in recidivism prediction, being marked as low-risk for committing a new crime is indisputably positive for the individual. On the contrary, a high activity goal -even though recommended- might not be realistic and thus advantageous for all individuals. Specifically, according to the Goal-Setting Theory by Latham and Locke [65], if an individual does not believe they can achieve their goal, they are unlikely to do so. Thus, users' goals should be close to their current abilities to hold sufficient motivational power.

To this end, we initially look into definitions based on predicted and actual outcomes, namely False Omission Rate (FOR), False Negative Rate (FNR), False Positive Rate (FPR) ratios, and Error Rate Ratio (ERR), that focus on erroneous predictions rather than solely positive outcomes (for definitions, formulas and interpretation, see Appendix B). However, we quickly notice that such metrics are prone to data biases and imbalances, as shown in the example below. Assume you have an imbalanced dataset of 3000 women -2000 low active and 1000 highly active- and 8000 men -3500 low active and 4500 highly active-. Imagine a model that misclassifies 100% of highly active women as low active, i.e., $FN = 1000$. Then, $ER_{women} = \frac{FP+FN}{P+N} = \frac{1000}{3000} = 33\%$. For men to have the same error rate, a model needs to misclassify 2640 highly active men ($FN = 2640$) as $ER_{men} = \frac{2640}{8000} = 33\%$. So, even though we have misclassified 100% of the highly active part of the minority group, by misclassifying only 59% ($\frac{2640}{8000}$) of the majority group, we can achieve in paper demographic parity with an ERR of 1.0 (optimal value). As discussed in Section 3, our data suffer from various biases and imbalances, and hence error-centric metrics would

not offer an objective comparison. Hence, for the purpose of this work, we utilize the widely used DIR, which is the ratio of base or selection rates between unprivileged and privileged groups, assuming equal ability across demographics:

$$\text{Disparate Impact Ratio} = \frac{\Pr(y^+ \mid G0)}{\Pr(y^+ \mid G1)}$$

where $y^+$ is the actual or predicted positive outcome label (base or selection rate, respectively), $G0$ is the minority (protected) group, and $G1$ is the majority group. Values less than 1 indicate that the majority group has a higher proportion of predicted positive outcomes than the minority group. A value of 1 indicates demographic parity. Values greater than 1 indicate that the minority group has a higher proportion of predicted positive outcomes than the majority one. For example, a value of 0.8 for a dataset with men/women as the majority/minority groups means that for every man receiving a high activity goal, only 0.8 women do so. According to the AIF360 toolkit (https://aif360.mybluemix.net/), accepted values are within [0.8,1.25], but such ranges are not universally accepted and might be adjusted on a task-by-task basis [29].

Having established our metric of choice, we move forward to our analysis of bias in the data generation (Section 3) and the model building and implementation (Section 4) streams.

## 3 EXPLORING BIAS IN PERSONAL INFORMATICS DATA GENERATION

Bias in the data generation stream can take the form of historical, representation, and measurement biases, as seen in Figure 1a. In this section, we explore all three sources, providing an answer to RQ2: *Are PI data susceptible to biases?*

### 3.1 Historical Bias

While historical biases cannot be measured directly in the specific dataset, there is evidence that PI is susceptible to several pre-existing or present biases. For completeness, we state the main findings of related literature below.

***Physical Activity Inequalities****.* In PI, physical activity data, such as step counts, are among the most common digital behavioral biomarkers. Similarly, in the MyHeart Counts dataset, they constitute the majority of the extracted HealthKit data. Specifically, the dataset includes 4920 users of step tracking compared to 626 users of sleep tracking, in line with previous research supporting that many users report not wearing their watch while sleeping [55]. However, inequalities in physical activity are well-reported [4, 50, 78]. Althoff et al. [4] use smartphone mobility data from over 68 million activity days by more than 700K individuals across 111 countries to quantify activity inequality. Their findings reveal variability in physical activity worldwide (measured in average step counts), where reduced activity in females explains a large portion of the observed activity inequality. Similarly, Guthold et al. [50] report that physical inactivity is twice as prevalent in high-income countries compared to low-income countries and they confirm lower activity levels in women than in men. Overall, the World Health Organization reports that "girls, women, older adults, people of low socioeconomic position, people with disabilities and chronic diseases, marginalized populations, indigenous people and the inhabitants of rural communities often have less access to safe, accessible, affordable and appropriate spaces and places in which to be physically active" [78]. Such inequalities, present in the real world, can undoubtedly creep into the behavioral data we build our models on.

***The Digital Divide****.* Similarly, as the world rapidly digitalizes, it threatens to exclude those that remain offline. Almost half the world's population, the majority of them women or citizens of developing countries, are still disconnected [74]. Even in the connected world, male internet users outnumber their female counterparts across regions. This "digital divide" encompasses even more discrepancies, such as the digital infrastructure quality and connectivity speed in rural or remote areas and the required skills to navigate technology [26]. Thus, it is

evident that data collected from any technological system, including PI, do not capture the entirety of the world population due to pre-existing inequalities in digital access and literacy.

**BYOD Study Design Biases**. On top of that, PI technologies are attracting attention as novel tools for data collection in clinical research, resulting in newfound demographic imbalances. Studies adopting a bring-your-own-device (BYOD) design, such as MyHeart Counts, are gaining traction because they are more user-friendly (participants use technologies they are already familiar with), achieve better participant compliance, potentially reduce the bias of introducing new technologies, and accelerate data collection from larger cohorts [27, 31]. However, the BYOD design may not support unbiased data collection from the target population where such technologies are intended to be deployed. In their work, Cho et al. [27] identify significant demographic disparities regarding race (50-85% white cohorts) in BYOD studies. Their findings align with the reported demographic divide existent in the composition of wearable users. Even though the gap is narrowing, a report by Ericsson ConsumerLab [38] documents that the majority of existing users of wearables are fit adults between 25–34 and that whilst females are more likely to own activity trackers, 63% of smartwatch owners are male. Hence, the technology used and the available participant cohorts in PI, especially for studies with BYOD design, such as the one under inspection, subject datasets to the same bias that has been exposed in the activity inequality and the digital divide literature.

## 3.2 Representation Bias

We discuss representation biases across three dimensions: misrepresented, underrepresented, and unevenly sampled populations.

**Misrepresented Populations**. Representation bias can emerge when the sample population does not reflect the general population (bias in rows). To evaluate for such biases in the MyHeart Counts dataset, we compare the ratios of majority and minority user segments as defined in Table 1 with the real-world ratios extracted from United States population censuses as the MyHeart Counts recruitment was spread across the country. Specifically, we utilize the United States Census Bureau (gender, race, and age [23] distributions), the Centers for Disease Control and Prevention (BMI [45, 46], joint issues [97], hypertension [42], and diabetes [77] distributions), and the American Heart Association (heart condition [98] distribution) data to extract the real distributions. Figure 2 showcases the results of this comparison in a radar plot. For example, while in the general United States population, we have approximately 1 female per 1 male (ratio of 1.0 in pink), in the MyHeart Counts HealthKit data, we have 0.2 females per 1 male, highlighting the strong underrepresentation of women in the dataset. The same applies to race, age, and hypertension segments, where the minority classes in the dataset (non-white users, users less than 45, and users with hypertension, respectively) do not reflect real-world ratios. An interesting finding is that, while in the United States, there exist approximately 0.3 underweight, overweight, or obese people for every person with normal weight, in the dataset, this ratio is doubled, in line with the research on BYOD design biases discussed above. Hence, potentially due to historical biases and study design choices, our analysis of the MyHeart Counts data (Figure 2) provides evidence that PI datasets might not represent the real target population.

**Underrepresented Populations**. PI datasets can still include underrepresented groups (bias in rows) even if sampled perfectly. Figure 4 shows significant imbalances, measured in the number of samples in the dataset, between minority and majority user segments across almost all protected attributes. We notice that even for representative sampling, e.g., users with joint or heart problems, the minority group is still significantly underrepresented in the data. Thus the model will likely be less robust for those few users with these conditions because it has fewer data to learn from. Overall, we see that the MyHeart Counts HealthKit data are skewed towards *white, fit males*, which needs to be considered in the preprocessing and model-building phases for a
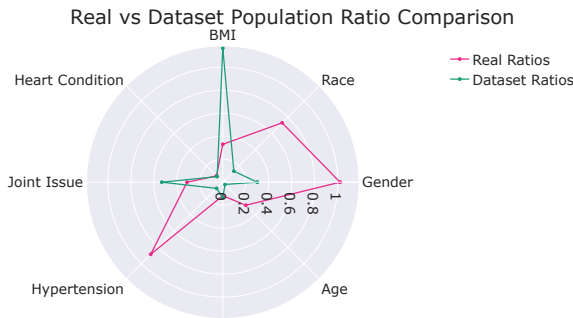
Fig. 2. Real (pink) versus dataset (green) ratios for different population segments in the MyHeart Counts dataset. The ratio is calculated as the number of the minority class divided by the number of the majority class instances. Larger distances between the two lines indicate larger deviations from the real ratios. We notice that attributes such as gender, age, race, and, to a smaller extent, hypertension and BMI suffer from representation bias.
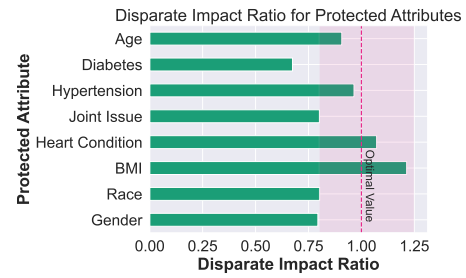


Fig. 3. A bar plot showing the DIR (ratio of base rates) per protected attribute. Values below 0.8 indicate bias against the minority segment, while values above 1.25 indicate bias against the majority segment. We notice that there exist biases in columns for diabetes patients, users with joint issues, and non-white minorities. While the data is borderline biased against women and people with non-healthy weight (underweight, overweight, or obese).

fairer machine learning life cycle. Note here that we cannot achieve realistic and equal representation unless the population is equally distributed. Ideally, a PI dataset should be representative of the target population but also large enough to consist of sufficient minority samples. In practice, this is challenging to achieve due to the effort and cost required to build large-scale PI datasets.
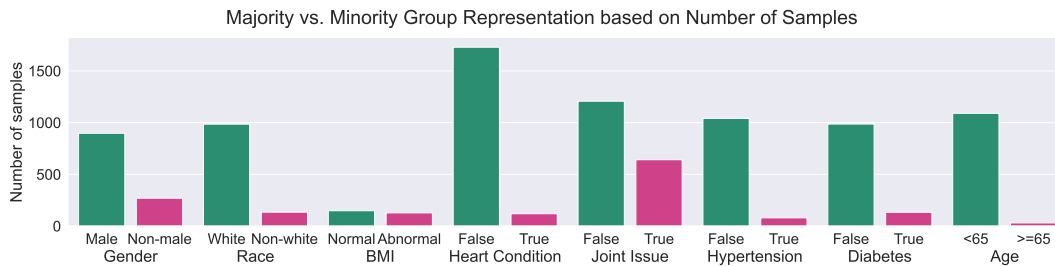


Fig. 4. A bar plot showcasing the number of samples per user segment split based on various protected attributes. We see significant underrepresentation of minority user segments across almost all attributes.

*Unevenly Sampled Populations*. Even if sampling is representative and equal (e.g., 50% male and 50% female users), the dataset can still suffer from representation bias if the sampling method is limited or uneven, e.g., all the males in the sample are highly active, but all females happen to be low active (bias in columns). This is also the case in the MyHeart Counts HealthKit data, as seen in Figure 3. The figure shows the DIR value (ratio of base rates, i.e., ratio of recorded high number of steps for unprivileged versus privileged groups) per protected attribute. A low value ($DIR < 0.8$) indicates a bias against the minority group, and a high value ($DIR > 1.25$) indicates a bias against the majority group. For our use case, a value of $DIR < 0.8$ means that the sample of the minority group is significantly less active than the sample of the majority group. For example, in the MyHeart

Counts data, diabetes patients, users with joint issues, racial minorities, and to a smaller extent, women, racial minorities, and overweight and obese users systematically perform lower step counts in the dataset compared to their majority segment counterparts. On the contrary, users of different age groups with or without hypertension or heart issues do not differ significantly in terms of step counts in the data.

### 3.3   Measurement Bias

In terms of measurement bias, we focus on the input modalities and their accuracy and discrepancies during data collection, i.e., we discuss how, in the MyHeart Counts data, the measurement method and accuracy vary across groups.

*Device Differences*. In the MyHeart Counts HealthKit dataset, data originate from different sources. Specifically, 33% comes from an iPhone, 11% comes from an Apple Watch, and 56% comes from multiple third parties. iPhones detect and calculate step counts through integrated sensors, such as an accelerometer, gyroscope, GPS, and in some models, a magnetometer. These sensor data are then analyzed by the motion coprocessor unit, namely a low-power unit that reads the sensors' output and makes the data available to applications via Apple's CoreMotion programming interface [11]. Specifically, it communicates with the CMMotionActivityManager [12], which is responsible for classifying whether the user is walking, running, in a vehicle, or stationary for periods of time. However, this process cannot be fully replicated in Apple watches due to inherent differences in placement (pocket versus wrist, fit, and usage habits. For instance, phones are known to underestimate user step counts due to non-carrying time in free-living conditions [5, 34]. On the contrary, Apple watches have been tested to be more accurate for measuring daily step counts for healthy adults [102].

Moreover, in the MyHeart Counts HealthKit data, there is also a statistically significant difference ($p < 0.05$) across segments: in Apple Watch ownership based on gender (46% of male participants have at least one watch entry compared to 28% non-males), heart condition (38% of participants with heart condition compared to 26% without), and ethnicity (41% of non-white participants compared to 36% white).

*Model Differences*. To make things worse, accuracy differences have been reported across consecutive generations of iPhone devices [34]. Incremental hardware changes may increase the quantity, modality, and quality of data available for the device to calculate the CMMotionActivityManager variables, which may improve the accuracy of activity recognition. For instance, iPhone 5S has introduced the M7 coprocessor; iPhones 6 and 6 Plus contain an M8 coprocessor; and 6S, 6S Plus, and SE have an M9 coprocessor, while prior models do not incorporate such a unit. The M8 has introduced the ability to differentiate between different activities [9], and the M9 has introduced "always-on" capabilities [10]. Additionally, newer versions of iOS may provide revised algorithms that improve recognition accuracy. In the MyHeartCounts HealthKit data, we encounter various iPhone models, starting from iPhone 4S (no coprocessor) and reaching iPhone 6S Plus (M9 coprocessor). We analyze whether differences in demographics correlate with differences in phone ownership, and we identify statistically significant differences ($p < 0.05$) based on gender and BMI. Specifically, females and people with normal BMI tend to own older and cheaper phones with fewer capabilities (see Figure 5).

*General Input Modality Differences*. Finally, most of the MyHeart Counts data comes from third parties, such as alternative wearables that communicate with the Apple Health app or fitness and well-being apps downloaded from the App Store. This is not uncommon in the PI domain, given the abundance and heterogeneity of available data sources. In our use case, we see, beyond the Apple Watch, Garmin, Polar, and Basis Peak wearable products, as well as various apps. With regards to third parties usage across segments, we identify statistically significant differences based on gender (91% of male participants have at least one third-party entry compared to 85% of non-male ones) and diabetes condition (97% of participants with diabetes have at least one third-party entry
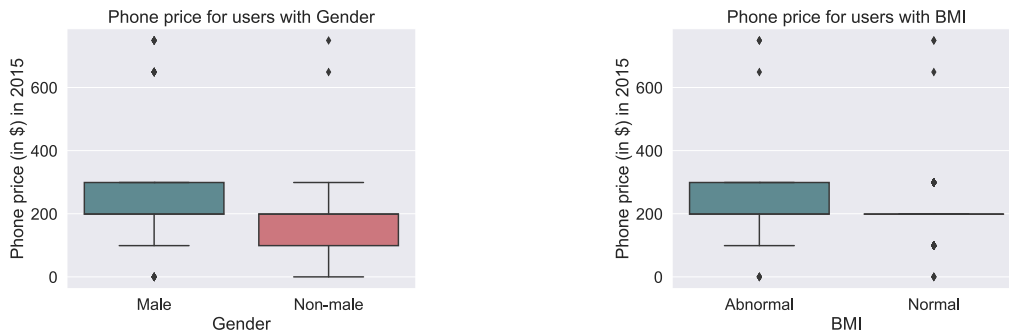
Fig. 5. Differences in the price of participants' phones as of September 2016 based on gender (left) and BMI (right). Females and people with BMI within the normal range tend to own older and cheaper phones with fewer capabilities.

compared to 90% without). However, different input devices or apps are proven to have different accuracies, likely to create measurement accuracy discrepancies between different users [36].

**Summary of biases in data generation:**
- Pre-existing *historical biases* are also present in digital biomarkers extracted from PI systems, due to well-documented phenomena, such as the global inequality in physical activity and the digital divide, leading to data generation that is not representative of the general population, which is also the case in our MyHeart Counts use case, where female, non-white, underweight, overweight or obese, young, and hypertensive users, are undersampled in the data.
- Even within well-sampled user groups, data imbalances, either in terms of user attributes or measured behaviors, are still prevalent due to realistic differences across user segments. Specifically, in our PI use case, we see significant underrepresentation of minority groups across all protected attributes and measured behavioral differences -not necessarily realistic- for users with diabetes, joint issues, non-healthy BMI, non-white users, and females.
- PI is susceptible to *measurement biases*, due to the heterogeneity in input modalities (smartphone versus smartwatch), performance and hardware differences across generations of devices, and usage of third-multiple party apps of unknown accuracy. Females are especially affected by such biases in our dataset, as they tend to own older devices with fewer capabilities and use to a greater extend multiple fitness-related third-party apps.

Given the awareness of certain historical, representation, and measurement biases in the data, practitioners can make informed decisions concerning appropriate preprocessing actions to alleviate potential negative effects. Such actions may include oversampling minority or undersampling majority user segments for misrepresented or underrepresented populations, choosing the appropriate sampling strategy to balance for unevenly sampled populations, or accounting for measurement differences across different devices or models.

## 4 EXPLORING BIAS IN PERSONAL INFORMATICS MODEL BUILDING AND IMPLEMENTATION

Bias in the model building and implementation stream can take the form of aggregation, learning, evaluation, and deployment biases, as seen in Figure 1b. In this section, we discuss all four sources, providing an answer to RQ3, namely: *Do machine learning models inherit PI data biases? Do they mitigate, propagate, or maybe even amplify them?*

## 4.1 Aggregation Bias

We evaluate aggregation bias by plotting the DIR (selection rate, i.e., rate of high activity goals predictions) for different user segments' predictions based on heart conditions, hypertension, joint issues, diabetes, race, BMI, gender, and age. Figure 6 shows the DIR scores for the segments, comparing data and baseline deep learning models. Specifically, we utilize two baseline models to capture the notions of "fairness through awareness" [35] and "fairness through unawareness" [63]. In fairness through awareness, fairness is captured by the principle that similar individuals should have similar classification outcomes. In our use case, the similarity is defined based on user demographics in the absence of other features. In practical terms, the aware model is trained on a feature set that includes protected attributes per user. On the other hand, fairness through unawareness is satisfied if no sensitive attributes are explicitly used in the learning process [103], namely, the unaware model is trained with features excluding protected attributes.

***Models' Description***. Our baseline models are sourced from prior work in the field of intelligent physical activity prediction, where Bampakis et al. [19], utilizing the MyHeart Counts dataset, benchmarked and evaluated six distinct learning paradigms from traditional machine learning models to advanced deep learning architectures. Their best model, a Long Short-Term Memory (LSTM) recurrent neural network, achieved a Mean Absolute Error (MAE) of 1087 steps, beating previous state-of-the-art approaches by 67% on the task of physical activity prediction.

We consider the following setting: we are given a time-series dataset $S = \{S^{G,0}, S^{G,1}\}$ of users segmented into two groups, $G0$ and $G1$, based on protected attribute $G$ (e.g., gender, age, etc.). The user data within each group are denoted as $S^{G,g} = \{s_1^{G,g}, \ldots, s_K^{G,g}\}$, where $g = \{0, 1\}$ and K is the number of users per group, conditioned on protected attribute $G$. Furthermore, the data of each user are stored as $s_i = \{X_i, y_i\}$, where input time series (step count values) of users $i = 1, \ldots, K$, are stored in $X_i \in \mathbb{R}^{D_x \times 1}$, where $D_x = 48$ (unaware model) is the length (in time steps) of a sample daily activity in the data, or $D_x = 56$ (aware model) is the length of a sample daily activity in the data plus the protected attribute features. Formally, our deep neural network architecture receives as input the users' daily activity samples ($X$) and passes them through LSTM layers with parameters $\theta_l = \{W_l, b_l\}$, weight matrix, and bias, respectively, for each layer $l$, to produce the output $\hat{y}$.

The optimization of the network parameters for LSTM layers is obtained by minimizing the binary cross entropy loss $\alpha_c$ defined as:

$$\Omega^* = \underset{\Omega = \{\theta_1, \ldots, \theta_3\}}{\arg\min} \alpha_c(\hat{y}, y) = \underset{\Omega = \{\theta_1, \ldots, \theta_3\}}{\arg\min} -\frac{1}{N} \sum_{i=1}^{N} \Big( y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \Big)$$

where $N$ represents the number of training samples *from both datasets* $\{S^{G,0}, S^{G,1}\}$.

We implement the proposed architectures in PyTorch Lightning [39]. The hyperparameter tuning is performed using the standard back-propagation algorithm and Adam optimizer with the default parameters [62]. To avoid overfitting in the deep models, we applied dropout with a varying portion of dropping nodes.

***Single Attribute Biases***. Our findings concerning machine learning model biases measured via DIR, as shown in Figure 6, highlight the following:
(1) Aware learning models are not foolproof against data biases in most cases (joint issues, diabetes, gender), and even amplify them for certain protected attributes (hypertension).
(2) Even excluding protected attributes from the training process of unaware models does not guarantee unbiased results in line with prior work [81]. Specifically, fairness through unawareness is also ineffective due to the presence of proxy features, namely attributes that work like proxies for protected attributes. Through such features, bias propagates from the data to models: for example, a person's walking behavior (measured in
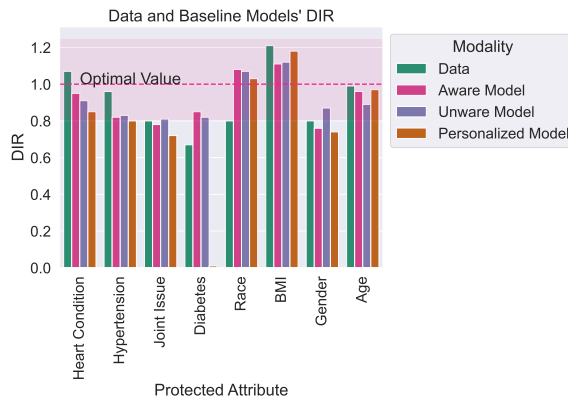
Fig. 6. A comparison of DIR between data, baseline model with protected attributes in the feature set (aware), and baseline model without protected attributes in the feature set (unaware). We see that the "one-size-fits-all" models propagate or, in some cases, amplify existing representation biases.
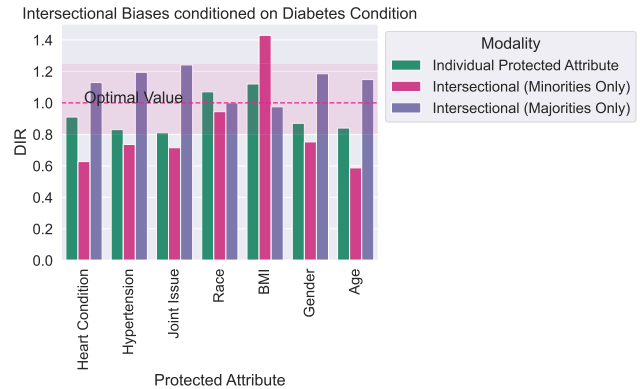
Fig. 7. A comparison of DIR given the unaware baseline model between user groups defined by a single protected attribute, e.g., gender, versus intersectional user groups defined by two attributes, e.g., gender and diabetes. Intersectional groups are either drawn from the minority or the majority classes for each attribute. The "one-size-fits-all" models' amplified biases are even more prevalent in intersectional cases.

step counts) is a good predictor of a person's gender, BMI, and age, which can thus be inferred, despite being hidden during training [60].
(3) Overall, diabetes patients have the largest bias gap compared to their non-diabetic counterparts, partially attributed to their highly biased training data to start with. Yet, users with hypertension have the largest difference between data and model biases since models trained on seemingly unbiased introduce bias during the learning process.

***Intersectional Biases****.* We also examine intersectional biases, as shown in Figure 7; namely we quantify the biases of the unaware model not only conditioned on an single protected attribute, but also on protected attribute combinations. Specifically, we consider two attributes at a time, and two different combination strategies: *minority-minority vs. rest* (e.g., diabetic women) and *majority-majority vs. rest* (e.g., non-diabetic men). Our results, which we present indicatively keeping the diabetes attribute fixed, highlight the widening intersectional biases for people who belong to more than one minority (in pink) across almost all attributes (with an exception of BMI, where people with non-healthy BMI are the majority group, despite usually being considered unprivileged in practice). The largest gap appears in people with more than one health condition, such as diabetic heart patients, and diabetic patients aged 65+. At the same time, people who do not belong to any minority groups (in purple), benefit across all attributes.

The trends in aggregation bias indicate that PI models do not tackle diverse user segments equally well, and reflect or even amplify representation biases existing in the data, especially when it comes to intersectional biases.

## 4.2 Learning Bias

In the PI literature, there has been a move toward personalization, straying from the "one-size-fits-all" mentality and its shortcomings, as discussed above. Contrary to generic models, personalized models are fine-tuned given the data of a single user or user segment. Accounting for such interindividual variability has been proven to

dramatically improve prediction performance in various tasks within the PI domain, such as pain detection, engagement estimation, and stress prediction from ubiquitous devices data [69, 89, 95]. Given the increasing popularity of the personalization paradigm, in this study, we investigate whether personalization as a modeling choice can amplify performance disparities across different user segments in the data, given the existence of representation bias.
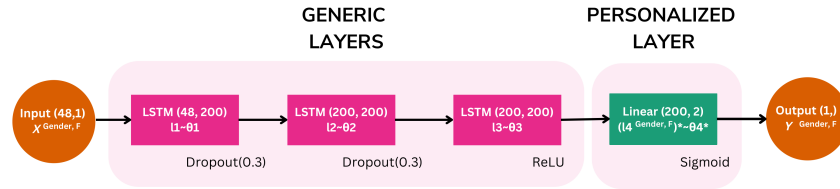


Fig. 8. Our personalized deep learning architecture inspired by CultureNet [87]. The last layer is indicatively fine-tuned based on gender for female users.

*Model Description.* We base our approach on the work of Rudovic et al. [87] and the CultureNet package [86] for building generalized and culturalized deep models to estimate engagement levels from face images of children with Autism Spectrum Condition. Specifically, we utilize our deep LSTM model, which is trained on data from all users, we freeze the network parameters $\{\theta_1, \ldots, \theta_3\}$ tuned to both minority and majority user groups, as described in Section 4.1, and then fine-tune the last layer ($\theta_4$), i.e., a linear fully-connected layer, to each user group separately based on the MyHeart Counts protected attributes (health condition, hypertension, joint issues, diabetes, race, BMI, gender, age). Figure 8 delineates the personalization process.

Formally, the learning during the fine-tuning process is attained through the last layer in the network, one for the minority and one for the majority user group. Before further optimization, the group-specific layers are initialized as $\theta_4^{G,0} \leftarrow \theta_4$ and $\theta_4^{G,1} \leftarrow \theta_4$, and then fine-tuned using the data from $G0$ ($S^{G,0}$) and $G1$ ($S^{G,1}$), respectively, *for each protected attribute G* as:

$$\left(\theta_4^{G,c}\right)^* = \arg\min_{\theta_4} -\frac{1}{N} \sum_{i=1}^{N \in S^{G,c}} \left(y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)\right), \quad c = \{0, 1\} \text{ and}$$

$$G = \{\text{gender, ethnicity, age, bmi, heart condition, hypertension, joint problem, diabetes}\}$$

The final network weights, $\theta_l = \{W_l, b_l\}$, are then used to perform the group-specific inference of next-day physical activity level from past behavior per protected attribute.

*Single Attribute Biases.* While we could not identify significant performance benefits either for the privileged or the unprivileged group by utilizing personalization in our use case, we encountered significant bias shortcomings of the approach. Specifically, across all protected attributes (with a borderline exception of race), we see that personalized models are more biased compared to either aware or unaware models or both. An extreme case appears in users with diabetes, where the personalized model "learns" that this user segment is less active than their non-diabetic counterpart in the dataset and thus provides them only with low activity goals, regardless of individual differences in physical activity levels. The intuition behind this behavior is that a personalized model is fine-tuned to a specific user segment, e.g., users with diabetes. If this segment suffers from representation bias in the dataset, which is true in many cases in PI, then personalized models amplify this bias through the fine-tuning process, as is evident in Figure 6. Our findings highlight that a common modeling choice in PI, such

as personalization, can negatively affect biases and asks for bias-aware personalization approaches to rip the benefits of user tailoring without leading to biased results.

## 4.3  Evaluation Bias

***Benchmark Selection***. In machine learning, models are optimized on their training data, but their quality is often evaluated based on benchmarks, such as ImageNet [32] in the computer vision community and MovieLens [51] in the recommender systems community. However, the ubiquitous computing community still suffers from a lack of benchmarks, or benchmarks that are limited to traditional tasks, such as human activity recognition [6] and sleep classification [25, 111]. To make things worse, oftentimes, benchmarks within the community are not representative of the target population. For example, within the fall detection domain, datasets usually comprise imitated falls performed by younger people while they are deployed on older people [93].
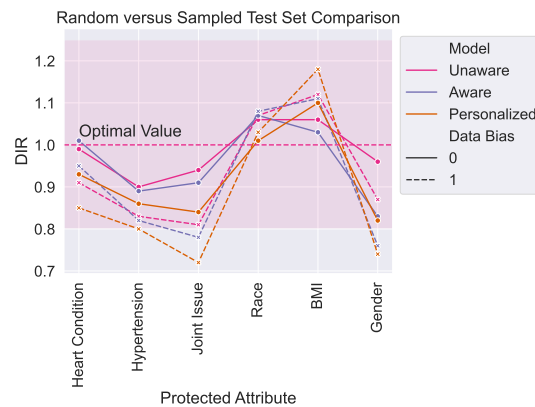


Fig. 9. A comparison of DIR between different test sets across models. We see that "perfect" test sets in terms of data bias (continuous lines) tend to hide imperfections in the trained models compared to the original test sets (dashed lines).

Yet, a misrepresentative benchmark encourages the development and deployment of models that perform well only on the data subset represented by the benchmark. To illustrate our point, given the lack of established benchmarks for our use case, we devise two distinct test sets for comparison purposes: our original (random) test set, $T1$, and a sampled subset of $T1$, $T0$, with demographic parity at base rate (DIR = 1.0). We then evaluate our models, namely the baseline aware and unaware and personalized LSTMs, on these two test sets. Figure 9 presents the results of our experimentation, where it is clear that $T0$, imitating a "perfect", fair world, consistently shows better performance concerning DIR compared to $T1$. Better performance is defined as smaller deviations from the optimal DIR value of 1.0. Essentially, an ideal-world benchmark, such as $T0$, is "hiding" the imperfections of our trained model, which has been proven to propagate or even amplify biases based on $T1$.

***Evaluation Metric Selection***. On a different note, evaluation bias can also emerge from the choice of metric used to quantify the models' performance. For instance, group fairness hybrid metrics, such as error rates, are prone to imbalances, as discussed earlier, and can hide disparities in other types of bias metrics, such as WAE metrics (see Appendix B). Similarly, aggregate measures, such as accuracy, can hide subgroup under-performance or conceal shortcomings in more important metrics for certain use cases, such as false positive or false negative rate [94].

## 4.4 Deployment Bias

***Changing Deployment Scenarios****.* We see at least two sources of deployment biases in PI. The first is related to the fact that the most active research areas within PI are Human-Activity Recognition and Sleep Classification. From this lens, FPs and FNs (Type I and Type II errors, respectively) in these scenarios are not critical, and models have been developed to maximize TPs. This dominant but limited view promotes deployment bias in novel use cases with the emergence of health-related intelligence embedded into PI systems. For example, given the novel ECG sensor data and AFib detection functionality, Type II errors should be minimized to avoid loss of life. It is thus critical to reassess the conceptualization of PI systems' evaluation practices and datasets and tailor them to their context.

***Development in Isolation****.* Second, learning models for PI systems are built and evaluated as if they were fully autonomous, while in reality, they operate in a complex socio-ethical system moderated by institutions and human decision-makers, also known as the "framing trap" [88]. Users may share their mHeatlh data with physicians for interpretation and disease management. Despite good performance in isolation, they may lead to harmful consequences because of human biases, such as confirmation bias. Specifically, physicians are more likely to believe AI that supports current practices, and opinions [79]. At the same time, research shows that physicians' perceptions about black male patients' physical activity behavior were significant predictors of their recommendations for coronary artery bypass graft surgery, independent of clinical factors, appropriateness, payer, and physician characteristics [100]. Such complicated interconnections highlight how evaluating a system in isolation creates unrealistic notions of its benefits and harms.

**Summary of biases in model building and implementation:**
- Digital biomarkers representation biases are propagated or even amplified by machine (deep) learning models, regardless of the inclusion of protected attributes in the feature set, due to the existence of proxy variables in PI data, e.g., steps, calories, that can be used by the model to infer hidden protected attributes. Such *aggregation biases* are also prevalent in our use case for users with joint issues, diabetes, hypertension, and female users.
- Common learning choices in PI, such as personalization, can introduce *learning biases*, if trained on biased data. In extreme cases, as highlighted in our use case for diabetic users, they can even introduce maximum bias, i.e., DIR = 0, while performing worse -in terms of bias- across all attributes.
- Our empirical results illustrate that model performance is highly susceptible to the representativeness of the PI benchmark used and highlight how *evaluation biases* can affect ubiquitous models in the evaluation phase.
- The application of machine learning in PI is not free of *deployment biases*, which can emerge from outdated evaluation practices emerging from the PI systems' early applications or the false assumption of autonomous PI systems' existence.

Given the awareness of certain aggregation and learning biases in the data, PI practitioners can make informed decisions concerning the machine learning paradigms or models to utilize or the necessary in-processing bias mitigation steps to apply. Additionally, aware of the evaluation biases present in PI data, they might choose to experiment with more benchmark datasets, evaluate their suitability for their target population apriori and select appropriate accuracy and fairness metrics to quantify performance across different user segments. Finally, they may follow a user-in-the-loop concept during the development phase, acknowledging the dependencies between PI systems and their users.

## 5 GENERALIZABILITY

This section aims to (i) demonstrate the straightforward applicability of our methodology and our open-source code [8] to other datasets and (ii) reveal initial insights about the generality of our findings and future steps.

While our analysis was conducted on the MyHeart Counts dataset, most of our findings can be generalized to other scenarios in PI and mHealth. To showcase this, we apply part of our experiments on two distinct datasets:

- **LifeSnaps:** LifeSnaps is a newly-released, multi-modal, time- and space-distributed dataset containing 71M rows of anthropological data, collected unobtrusively for the total course of more than four months by 71 participants. Based on data availability, we consider three protected attributes in Lifesnaps, namely gender, age, and BMI. Also, given the lack of official benchmark tasks, namely tasks built specifically on this dataset that are selected to be representative of relevant machine learning workloads and to evaluate competing models, we consider the "next-day physical activity prediction" task for model training, same with the MyHeart Counts dataset.
- **MIMIC-III:** MIMIC-III is an established, large-scale clinical dataset consisting of information concerning more than 38K patients admitted to intensive care units (ICU) at a large tertiary care hospital. Based on data availability, in MIMIC-III, we consider six protected attributes, namely gender, ethnicity, language, insurance, religion, and age. Contrary to LifeSnaps or MyHeart Counts, there exists a public benchmark suite that includes four different clinical prediction tasks for MIMIC-III [52]. For this analysis, we utilize the "in-hospital mortality" task as a binary classification equivalent to the "next-day physical activity prediction" task.
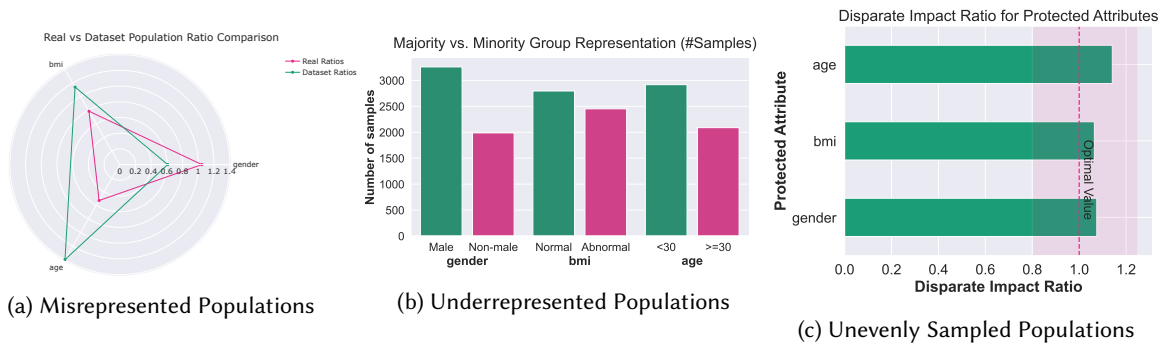


(a) Misrepresented Populations

(b) Underrepresented Populations

(c) Unevenly Sampled Populations

Fig. 10. LifeSnaps data representation biases



(a) Misrepresented Populations

(b) Underrepresented Populations
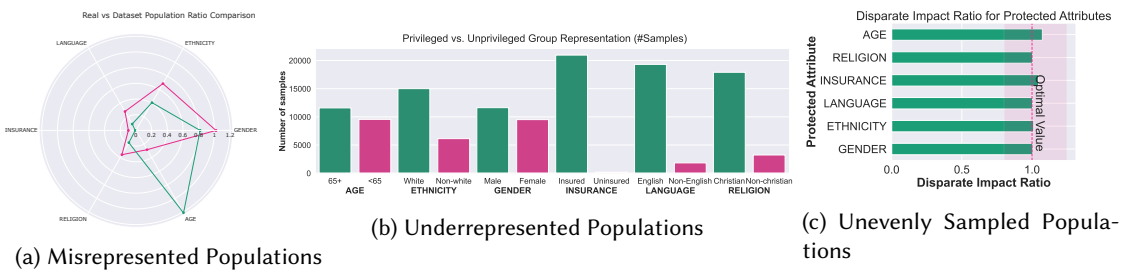
(c) Unevenly Sampled Populations

Fig. 11. MIMIC-III data representation biases

In exploring biases, we identified both commonalities and differences across PI datasets. Regarding the data generation stream, *representation biases seem to be the norm in PI datasets*, naturally leading to *learning and aggregation biases* in the model building and implementation stream and highlighting the need for increased awareness among researchers and practitioners in the field. Having said that, the identified biases are distinct in each dataset, emerging mostly from their recruitment methodology and the availability of protected attributes.

***Bias in Rows Commonalities***. All three datasets suffer from some type of "bias in rows" as seen in Figures 10a and 11a. Specifically, both LifeSnaps and MIMIC-III suffer from misrepresented populations. In LifeSnaps (Figure 10a), younger people are overrepresented due to university-based recruitment, while in MIMIC-III (Figure 11a) older people are overrepresented due to ICU-based recruitment. Additionally, while gender and ethnicity representation is improved compared to MyHeart Counts, still white males are overrepresented in all three datasets.

MIMIC-III, similarly to MyHeartCounts, suffers from underrepresented populations, such as uninsured, non-white, non-English-speaking, or non-christian users (Figure 11b). These biases are, in turn, propagated to the baseline learning models, in line with prior work [85].

***Bias in Columns Differences***. When "bias in columns" is explored, contrary to the MyHeart Counts data, both datasets are evenly sampled in terms of outcome labels, namely physical activity in LifeSnaps and in-hospital mortality for MIMIC-III (Figures 10c and 11c, respectively). Hence, these findings may not directly generalize to other PI datasets but are still included in our methodology for completeness and visibility. Specifically, contrary to population demographics which can capture misrepresented and underrepresented groups, an analysis for unevenly sampled populations is not commonly performed during data exploration, whereas in certain cases, such as MyHeart Counts, it could reveal behavioral discrepancies across populations.

Overall, these findings concerning generalizability highlight the need for comprehensive data and model evaluation in PI and, by extension, mHealth. It is high time PI researchers and practitioners looked beyond performance-only metrics to human-centric metrics capturing biases and demographic parity.

## 6 RELATED WORK

With the widespread adoption of intelligent systems and applications in our everyday lives, accounting for fairness has gained significant traction in designing and deploying systems. Specifically, fairness has been studied extensively in domains such as natural language processing [21, 43], recommender systems [66, 104, 106], and computer vision [22, 107]. Yet, evidence for fairness in the PI setting is lacking. Closer to PI, fairness research in the healthcare setting is still in its infancy [40]. The digitization of medical data has enabled the scientific community to collect large amounts of heterogeneous, multi-modal data and develop machine learning algorithms for a variety of medical tasks. During the process, various fairness limitations have been uncovered based on the three most prominent data types, namely medical image data, structured electronic health record (EHR) data, and textual data.

First, medical imaging has been the most widely used data source for machine learning in healthcare, and biases in them have received attention [56]. For example, Larrazabal et al. [64] utilize two commonly used X-ray image datasets to diagnose various chest diseases under different gender imbalance conditions and showcase that the minority gender group systematically performs worse than the majority gender group. Similarly, according to Adamson and Smith [1], relying on machine learning for skin cancer screening may exacerbate potential racial disparities in dermatology.

On a different note, EHR systems store multi-modal, heterogeneous patient data, such as demographics, diagnoses, and clinical records, and have been used for various tasks such as medical concept extraction, mortality prediction, and disease inference. Regarding EHR data fairness, Meng et al. [73] identify race-level differences in the predictions of neural network models on the MIMIC-IV dataset [57], with Black and Hispanic patients being less likely to receive interventions or receiving interventions of shorter average duration. Similarly, Röösli et al. [85] reveal a strong class imbalance problem and significant fairness concerns for Black and publicly insured ICU patients in the same dataset.

Concerning textual EHR data, Chen et al. [24] examine clinical and psychiatric notes to predict intensive care unit mortality and 30-day psychiatric readmission. Their analysis reveals differences in prediction accuracy,

and biases are present in terms of gender and insurance type for mortality prediction, and insurance policy for psychiatric 30-day readmission. Within the same scope, Zhang et al. [112] train deep embedding models on medical notes from the MIMIC-III database [58], and find that classifiers trained from their embeddings exhibit statistically significant differences in performance, often favoring the majority group regarding gender, language, ethnicity, and insurance status.

Yet, despite the emerging research on fairness in healthcare, its proximity to PI, and the widespread adoption of PI technologies, biases in PI have been barely explored. An initial effort of capturing biases in digital biomarkers is reported by Paviglianiti and Pasero [80]. Their Vital-ECG, a wearable smart device that collects electrocardiogram and plethysmogram signals, is embedded with machine learning algorithms to monitor arterial blood pressure and is found to underestimate the risk of disease in female patients. While this is a first step in uncovering biases in PI, it is far from a complete study of bias in PI. As highlighted by research in other domains, bias has multiple facets that may affect system fairness. To this end, our work aims to raise awareness and set up a systematic approach for a comprehensive analysis of data and machine learning model biases/fairness in PI systems.

## 7 DISCUSSION & CONCLUSIONS

This paper presents the first-of-its-kind, in-depth study of bias in PI by analyzing the most extensive digital biomarkers data to date. We provide empirical and analytical evidence of sources of bias at every stage of the PI ML development pipeline, from data ingestion to model deployment. In response to *RQ1*, we recognize the limitations of hybrid group fairness metrics in overcoming data imbalances and conclude that there is no optimal metric as of now capturing the idiosyncrasies of PI. Additionally, in response to *RQ2* and *RQ3*, we show that bias exists across all stages of the machine learning lifecycle, both in the data generation and model building and implementation streams. Different user minorities are affected by diverse types of bias, but users with diabetes, joint issues, or hypertension and female users show higher degrees of impact adversity in our MyHeart Counts use case due to representation, aggregation, and learning biases. Our findings echo concerns similar to those raised in the evaluation for healthcare technologies [3]. While some of our findings are specific to the investigated use case, they can, for the most part, be extended to PI tasks more broadly. Below we present limitations of our work that create new opportunities for future research and provide recommendations for future work for studying and mitigating bias in PI.

### 7.1 Limitations

*Alternative PI Use Cases.* Our work presents the first study of bias in PI research and development and does not study and compare bias in commercial PI systems, such as consumer smartphones and wearables, which we position as a future work direction. This is due to the prevalence of commercial black box models -which can be attributed to competitive advantage in an emerging market- and closed data because of ethical and privacy considerations. Yet, such restrictions have limited our use case, which may not seem as critical as AFib detection, for instance, but provides a strong indication of how PI data and models are susceptible to bias and is large enough to ensure the generalizability of our findings. Hence, our findings should be interpreted with these limitations in mind and not be seen as a generic evaluation of bias across all PI systems.

*In-the-wild Data Quantity versus Data Quality.* In our search for large-scale data, we had to partly sacrifice data quality (e.g., missing values, noise, duplicate measurements), as often happens with in-the-wild datasets of the scale of MyHeart Counts. Nevertheless, we engaged in thorough preprocessing methods benchmarking to ensure the best possible quality for our training data, as reported in prior work [19]. Due to the small sample sizes for certain user groups conditioned on a protected attribute, e.g., Pacific Islanders or American Indians for the ethnicity attribute, we had to binarize all protected attributes to avoid immense imbalances between majority and minority groups. Nevertheless, we recognize that some minority groups might be treated more unfairly than

others by the data and algorithms, a fact not captured in the current configuration. On the same note, gender was treated as a binary concept in the MyHeart Counts dataset, and recognizing diverse gender identities was outside of our control for the purpose of this study.

## 7.2 Future Work Directions

*Inclusive Training and Evaluation Datasets for Real-life Scenarios.* Appropriate PI datasets for fueling future bias research in the domain are still lacking. Due to the sensitivity of the data at hand, many datasets are proprietary with restrictive Institutional Review Board (IRB) agreements. Out of the open PI datasets, most are small-scale [110] due to the high effort and equipment cost required for building larger datasets or are conducted in-the-lab failing to represent the target population. To this end, any future work publishing open, large-scale, in-the-wild PI data sourced from diverse populations in terms of geographic location, gender, age, and health conditions, is significantly contributing to the advancement of the domain. Also, given the prevalence of small-scale datasets, future work should focus on quantifying biases in small digital biomarkers data, as realistically, most institutions will never acquire big data [17]. Additionally, due to the recentness of the domain and the closed-sourced data and algorithms, there is a lack of established benchmarks, especially regarding emerging PI tasks, such as fertility prediction, AFib, or fall detection. To this end, similarly to the work of Harutyunyan et al. [52], which published benchmarks for electronic health records tasks, future work should create inclusive and representative benchmarks for tasks within the PI domain.

*Fairness Metrics Capturing PI Idiosyncrasies.* Digital biomarkers are essentially sequential time-series data, inherently different from images and audio, where fairness research is most advanced. Hence, there is work to be done in quantifying bias and identifying idiosyncrasies in sequential physiological and behavioral data. For instance, many PI tasks are formulated as regression problems, but regression-specific fairness metrics are limited in the literature [49]. Beyond that, future work should explore diverse definitions of bias and their suitability for heterogeneous PI tasks. For example, how is demographic disparity defined in AFib detection, where false negatives can be fatal, versus human-activity recognition, where false positives deteriorate the user experience? Also, which metrics are most appropriate in tasks with no clear positive outcome, such as fertility prediction, given the shortcomings of error-based metrics as discussed in Section 2.3? The latter research gap provides opportunities for future work in redefining error-based fairness metrics that are more robust to representation biases in the data, as is the case in digital biomarkers.

*Benchmarking Bias Mitigation Approaches in PI.* While the focus of this work is uncovering the susceptibility of digital biomarkers to data and model biases, there is plenty of work to be done in benchmarking preprocessing, in-processing, and post-processing bias mitigation approaches or developing new ones for capturing the idiosyncrasies of digital biomarkers and their respective PI tasks. PI tasks are undoubtedly different from learning to rank or classification scenarios not only because of the nature of their data but also oftentimes their problem formulation as regression. Yet, fair regression is considerably overlooked compared to fair classification [2], and most fairness libraries (AIF360, FairLearn) feature only a few bias mitigation algorithms' implementations for regression tasks. Additionally, PI literature uses discrete machine learning paradigms, such as self-supervised learning, multi-task learning, and personalized learning, among others, whose consequences to algorithm bias are yet to be explored. Finally, due to privacy considerations for sensitive digital biomarkers, many times PI data are not accompanied by protected attributes for the population they describe, making it cumbersome to perform a fairness evaluation. To this end, future work should investigate the space of "fairness in unawareness", or, in other words, how you can quantify and mitigate biases in the absence of protected attributes.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Adewole S Adamson and Avery Smith. 2018. Machine learning and health care disparities in dermatology. *JAMA dermatology* 154, 11 (2018), 1247–1248.

[2] Alekh Agarwal, Miroslav Dudík, and Zhiwei Steven Wu. 2019. Fair regression: Quantitative definitions and reduction-based algorithms. In *International Conference on Machine Learning*. PMLR, JMLR, Campridge, MA, United States, 120–129.

[3] Muhammad Aurangzeb Ahmad, Arpit Patel, Carly Eckert, Vikas Kumar, and Ankur Teredesai. 2020. Fairness in machine learning for healthcare. In *Proceedings of the 26th acm sigkdd international conference on knowledge discovery & data mining*. ACM, New York, NY, United States, 3529–3530.

[4] Tim Althoff, Rok Sosič, Jennifer L Hicks, Abby C King, Scott L Delp, and Jure Leskovec. 2017. Large-scale physical activity data reveal worldwide activity inequality. *Nature* 547, 7663 (2017), 336–339.

[5] Shiho Amagasa, Masamitsu Kamada, Hiroyuki Sasai, Noritoshi Fukushima, Hiroyuki Kikuchi, I-Min Lee, Shigeru Inoue, et al. 2019. How well iPhones measure steps in free-living conditions: cross-sectional validation study. *JMIR mHealth and uHealth* 7, 1 (2019), e10418.

[6] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra Perez, and Jorge Luis Reyes Ortiz. 2013. A public domain dataset for human activity recognition using smartphones. In *Proceedings of the 21th international European symposium on artificial neural networks, computational intelligence and machine learning*. MDPI, Basel, Switzerland, 437–442.

[7] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias. In *Ethics of Data and Analytics*. Auerbach Publications, Boca Raton, FL, United States, 254–264.

[8] Anonymous. 2023. Source Code for "Uncovering Bias in Personal Informatics". [Accessed 4-Jan-2023].

[9] Apple. 2014. Apple — September Event 2014 — youtube.com. https://www.youtube.com/watch?v=38IqQpwPe7s&ab_channel=Apple. [Accessed 14-Nov-2022].

[10] Apple. 2015. Apple - September Event 2015 — youtube.com. https://www.youtube.com/watch?v=0qwALOOvUik&ab_channel=Apple. [Accessed 14-Nov-2022].

[11] Apple. 2022. Apple Developer Documentation — developer.apple.com. https://developer.apple.com/documentation/coremotion. [Accessed 14-Nov-2022].

[12] Apple. 2022. Apple Developer Documentation — developer.apple.com. https://developer.apple.com/documentation/coremotion/cmmotionactivitymanager. [Accessed 14-Nov-2022].

[13] Apple. 2022. Empowering people to live a healthier day. *Health Report* 1 (2022), 60.

[14] Apple. 2022. How Apple is empowering people with their health information — apple.com. https://www.apple.com/gr/newsroom/2022/07/how-apple-is-empowering-people-with-their-health-information/. [Accessed 14-Nov-2022].

[15] Apple. 2022. ResearchKit — researchkit.org. http://researchkit.org/. [Accessed 14-Nov-2022].

[16] Yazeed Awwad, Richard Fletcher, Daniel Frey, Amit Gandhi, Maryam Najafian, and Mike Teodorescu. 2020. *Exploring fairness in machine learning for international development*. Technical Report. CITE MIT D-Lab.

[17] Ricardo Baeza-Yates. 2018. BIG, small or Right Data: Which is the proper focus.

[18] Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. 2019. Differential privacy has disparate impact on model accuracy. *Advances in neural information processing systems* 32 (2019), 15479–-15488.

[19] Asterios Bampakis, Sofia Yfantidou, and Athena Vakali. 2022. UBIWEAR: An end-to-end, data-driven framework for intelligent physical activity prediction to empower mHealth interventions. In *2022 IEEE International Conference on E-health Networking, Application & Services (HealthCom)*. IEEE, IEEE, New York, NY, United States, 56–62.

[20] Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, Seema Nagar, Karthikeyan Natesan Ramamurthy, John Richards, Diptikalyan Saha, Prasanna Sattigeri, Moninder Singh, Kush R. Varshney, and Yunfeng Zhang. 2018. AI Fairness 360: An Extensible Toolkit for Detecting, Understanding, and Mitigating Unwanted Algorithmic Bias. https://arxiv.org/abs/1810.01943

[21] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems* 29 (2016), 4356–-4364.

[22] Joy Buolamwini and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*. PMLR, JMLR, Campridge, MA, United States, 77–91.

[23] United States Census Bureu. 2022. National Population by Characteristics: 2020-2021. https://www.census.gov/data/tables/time-series/demo/popest/2020s-national-detail.html

[24] Irene Y Chen, Peter Szolovits, and Marzyeh Ghassemi. 2019. Can AI help reduce disparities in general medical and mental health care? *AMA journal of ethics* 21, 2 (2019), 167–179.

[25] Xiaoli Chen, Rui Wang, Phyllis Zee, Pamela L Lutsey, Sogol Javaheri, Carmela Alcántara, Chandra L Jackson, Michelle A Williams, and Susan Redline. 2015. Racial/ethnic differences in sleep disturbances: the Multi-Ethnic Study of Atherosclerosis (MESA). *Sleep* 38, 6 (2015), 877–888.

[26] Krish Chetty, Liu Qigui, Nozibele Gcora, Jaya Josie, Li Wenwei, and Chen Fang. 2018. Bridging the digital divide: measuring digital literacy. *Economics* 12, 1 (2018), 1–17.

[27] Peter Jaeho Cho, Jaehan Yi, Ethan Ho, Md Mobashir Hasan Shandhi, Yen Dinh, Aneesh Patil, Leatrice Martin, Geetika Singh, Brinnae Bent, Geoffrey Ginsburg, et al. 2022. Demographic Imbalances Resulting From the Bring-Your-Own-Device Study Design. *JMIR mHealth and uHealth* 10, 4 (2022), e29510.

[28] Erin Collins. 2018. Punishing risk. *Geo. LJ* 107 (2018), 57.

[29] Equal Employment Opportunity Commission et al. 1990. Uniform guidelines on employee selection procedures. *Fed Register* 1 (1990), 216–243.

[30] Jeffrey Dastin. 2018. Amazon scraps secret AI recruiting tool that showed bias against women. In *Ethics of Data and Analytics*. Auerbach Publications, Boca Raton, FL, United States, 296–299.

[31] Charmaine Demanuele, Cynthia Lokker, Krishna Jhaveri, Pirinka Georgiev, Emre Sezgin, Cindy Geoghegan, Kelly H Zou, Elena Izmailova, and Marie McCarthy. 2022. Considerations for Conducting Bring Your Own "Device"(BYOD) Clinical Studies. *Digital biomarkers* 6, 2 (2022), 47–60.

[32] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. IEEE, IEEE, New York, NY, United States, 248–255.

[33] Emily Denton, Alex Hanna, Razvan Amironesei, Andrew Smart, Hilary Nicole, and Morgan Klaus Scheuerman. 2020. Bringing the people back in: Contesting benchmark machine learning datasets. *arXiv preprint arXiv:2007.07399* (2020), 1–6.

[34] Markus J Duncan, Kelly Wunderlich, Yingying Zhao, and Guy Faulkner. 2018. Walk this way: validity evidence of iphone health application step count in laboratory and free-living conditions. *Journal of sports sciences* 36, 15 (2018), 1695–1704.

[35] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. ACM, New York, NY, United States, 214–226.

[36] Fatema El-Amrawy and Mohamed Ismail Nounou. 2015. Are currently available wearable devices for activity tracking and heart rate monitoring accurate, precise, and medically beneficial? *Healthcare informatics research* 21, 4 (2015), 315–320.

[37] Daniel A Epstein, Clara Caldeira, Mayara Costa Figueiredo, Xi Lu, Lucas M Silva, Lucretia Williams, Jong Ho Lee, Qingyang Li, Simran Ahuja, Qiuer Chen, et al. 2020. Mapping and taking stock of the personal informatics literature. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 4 (2020), 1–38.

[38] Ericsson ConsumerLab. 2016. Wearable technology and the IoT.

[39] William Falcon et al. 2019. PyTorch Lightning.

[40] Qizhang Feng, Mengnan Du, Na Zou, and Xia Hu. 2022. Fair machine learning in healthcare: A review. *arXiv preprint arXiv:2206.14397* (2022), 1–21.

[41] Will Fleisher. 2021. What's Fair about Individual Fairness?. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, New York, NY, United States, 480–490.

[42] Centers for Disease Control, Prevention (CDC), et al. 2019. Hypertension cascade: hypertension prevalence, treatment and control estimates among US adults aged 18 years and older applying the criteria from the American College of Cardiology and American Heart Association's 2017 Hypertension Guideline—NHANES 2013–2016.

[43] William R Frey, Desmond U Patton, Michael B Gaskell, and Kyle A McGregor. 2020. Artificial intelligence and inclusion: Formerly gang-involved youth as domain experts for analyzing unstructured twitter data. *Social Science Computer Review* 38, 1 (2020), 42–56.

[44] Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2021. The (im) possibility of fairness: Different value systems require different mechanisms for fair decision making. *Commun. ACM* 64, 4 (2021), 136–143.

[45] Cheryl D Fryar, Margaret D Carroll, and Joseph Afful. 2020. Prevalence of underweight among adults aged 20 and over: United States, 1960–1962 through 2017–2018.

[46] Cheryl D Fryar, Margaret D Carroll, Joseph Afful, et al. 2020. Prevalence of overweight, obesity, and severe obesity among adults aged 20 and over: United States, 1960–1962 through 2017–2018. , 7 pages.

[47] Andrew Garbett, David Chatting, Gerard Wilkinson, Clement Lee, and Ahmed Kharrufa. 2018. ThinkActive: designing for pseudonymous activity tracking in the classroom. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, United States, 1–13.

[48] The Radicati Group. 2021. Mobile Statistics Report, 2021-2025. https://www.radicati.com/wp/wp-content/uploads/2021/Mobile_Statistics_Report,_2021-2025_Executive_Summary.pdf

[49] Furkan Gursoy and Ioannis A Kakadiaris. 2022. Error Parity Fairness: Testing for Group Fairness in Regression Tasks. *arXiv preprint arXiv:2208.08279* (2022), 1–12.

[50] Regina Guthold, Gretchen A Stevens, Leanne M Riley, and Fiona C Bull. 2018. Worldwide trends in insufficient physical activity from 2001 to 2016: a pooled analysis of 358 population-based surveys with 1· 9 million participants. *The lancet global health* 6, 10 (2018), e1077–e1086.

[51] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)* 5, 4 (2015), 1–19.

[52] Hrayr Harutyunyan, Hrant Khachatrian, David C Kale, Greg Ver Steeg, and Aram Galstyan. 2019. Multitask learning and benchmarking with clinical time series data. *Scientific data* 6, 1 (2019), 1–18.

[53] Steven G Hershman, Brian M Bot, Anna Shcherbina, Megan Doerr, Yasbanoo Moayedi, Aleksandra Pavlovic, Daryl Waggott, Mildred K Cho, Mary E Rosenberger, William L Haskell, et al. 2019. Physical activity, sleep and cardiovascular health data for 50,000 individuals from the MyHeart Counts Study. *Scientific data* 6, 1 (2019), 1–10.

[54] Kelly M Hoffman, Sophie Trawalter, Jordan R Axt, and M Norman Oliver. 2016. Racial bias in pain assessment and treatment recommendations, and false beliefs about biological differences between blacks and whites. *Proceedings of the National Academy of Sciences* 113, 16 (2016), 4296–4301.

[55] Hayeon Jeong, Heepyung Kim, Rihun Kim, Uichin Lee, and Yong Jeong. 2017. Smartwatch wearing behavior analysis: a longitudinal study. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 3 (2017), 1–31.

[56] Fei Jiang, Yong Jiang, Hui Zhi, Yi Dong, Hao Li, Sufeng Ma, Yilong Wang, Qiang Dong, Haipeng Shen, and Yongjun Wang. 2017. Artificial intelligence in healthcare: past, present and future. *Stroke and vascular neurology* 2, 4 (2017), 230–243.

[57] Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2020. Mimic-iv.

[58] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific data* 3, 1 (2016), 1–9.

[59] Joseph Jofish Kaye, Mary McCuistion, Rebecca Gulotta, and David A Shamma. 2014. Money talks: tracking personal finances. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, United States, 521–530.

[60] Andrei Kazlouski, Thomas Marchioro, and Evangelos Markatos. 2022. What your Fitbit Says about You: De-anonymizing Users in Lifelogging Datasets. In *Proceedings of the 18th International Conference on Security and Cryptography - SECRYPT*. SciTePress, Setubal, Portugal, 806–811.

[61] Young-Ho Kim, Jae Ho Jeon, Eun Kyoung Choe, Bongshin Lee, KwonHyun Kim, and Jinwook Seo. 2016. TimeAware: Leveraging framing effects to enhance personal productivity. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, United States, 272–283.

[62] Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *Proceedings of the International Conference on Learning Representations (ICLR)* (2015).

[63] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. *Advances in neural information processing systems* 30 (2017).

[64] Agostina J Larrazabal, Nicolás Nieto, Victoria Peterson, Diego H Milone, and Enzo Ferrante. 2020. Gender imbalance in medical imaging datasets produces biased classifiers for computer-aided diagnosis. *Proceedings of the National Academy of Sciences* 117, 23 (2020), 12592–12594.

[65] Gary P Latham and Edwin A Locke. 1991. Self-regulation through goal setting. *Organizational behavior and human decision processes* 50, 2 (1991), 212–247.

[66] Jurek Leonhardt, Avishek Anand, and Megha Khosla. 2018. User fairness in recommender systems. In *Companion Proceedings of the The Web Conference 2018*. ACM, New York, NY, United States, 101–102.

[67] Ian Li, Anind Dey, and Jodi Forlizzi. 2010. A stage-based model of personal informatics systems. In *Proceedings of the SIGCHI conference on human factors in computing systems*. ACM, New York, NY, United States, 557–566.

[68] Zhiguo Li, Subhro Das, James Codella, Tian Hao, Kun Lin, Chandramouli Maduri, and Ching-Hua Chen. 2019. An Adaptive, Data-Driven Personalized Advisor for Increasing Physical Activity. *IEEE journal of biomedical and health informatics* 23, 3 (May 2019), 999–1010. https://doi.org/10.1109/jbhi.2018.2879805

[69] Daniel Lopez-Martinez and Rosalind Picard. 2017. Multi-task neural networks for personalized pain recognition from physiological signals. In *2017 Seventh International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW)*. IEEE, IEEE, New York, NY, United States, 181–184.

[70] Steven A Lubitz, Anthony Z Faranesh, Caitlin Selvaggi, Steven J Atlas, David D McManus, Daniel E Singer, Sherry Pagoto, Michael V McConnell, Alexandros Pantelopoulos, and Andrea S Foulkes. 2022. Detection of atrial fibrillation in a large population using wearable devices: the Fitbit heart study. *Circulation* 146, 19 (2022), 1415–1424.

[71] Alexander Malafeev, Dmitry Laptev, Stefan Bauer, Ximena Omlin, Aleksandra Wierzbicka, Adam Wichniak, Wojciech Jernajczyk, Robert Riener, Joachim Buhmann, and Peter Achermann. 2018. Automatic human sleep stage scoring using deep neural networks. *Frontiers in neuroscience* 12 (2018), 781.

[72] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* 54, 6 (2021), 1–35.

[73] Chuizheng Meng, Loc Trinh, Nan Xu, and Yan Liu. 2021. Mimic-if: Interpretability and fairness evaluation of deep learning models on mimic-iv dataset.

[74] Amina Mohammed. 2021. With Almost Half of World's Population Still Offline, Digital Divide Risks Becoming 'New Face of Inequality,'Deputy Secretary-General Warns General Assembly.

[75] Deirdre K Mulligan, Joshua A Kroll, Nitin Kohli, and Richmond Y Wong. 2019. This thing called fairness: Disciplinary confusion realizing a value in technology. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–36.

[76] Arvind Narayanan. 21. Fairness definitions and their politics. In *Tutorial presented at the Conf. on Fairness, Accountability, and Transparency.*

[77] US Department of Health, Human Services, et al. 2020. National diabetes statistics report, 2020.

[78] World Health Organization. 2019. *Global action plan on physical activity 2018-2030: more active people for a healthier world.* World Health Organization, Geneva, Switzerland.

[79] Ravi B Parikh, Stephanie Teeple, and Amol S Navathe. 2019. Addressing bias in artificial intelligence in health care. *Jama* 322, 24 (2019), 2377–2378.

[80] Annunziata Paviglianiti and Eros Pasero. 2020. VITAL-ECG: a de-bias algorithm embedded in a gender-immune device. In *2020 IEEE International Workshop on Metrology for Industry 4.0 & IoT*. IEEE, IEEE, New York, NY, United States, 314–318.

[81] Dino Pedreshi, Salvatore Ruggieri, and Franco Turini. 2008. Discrimination-aware data mining. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, New York, NY, United States, 560–568.

[82] Dana Pessach and Erez Shmueli. 2022. A Review on Fairness in Machine Learning. *ACM Computing Surveys (CSUR)* 55, 3 (2022), 1–44.

[83] Kelly Peterson, Ognjen Rudovic, Ricardo Guerrero, and Rosalind W Picard. 2017. Personalized gaussian processes for future prediction of alzheimer's disease progression. In *Proceedings of the Machine Learning for Health Workshop - ML4H)*. ACM, New York, NY, United States, 1–13.

[84] Inioluwa Deborah Raji, Timnit Gebru, Margaret Mitchell, Joy Buolamwini, Joonseok Lee, and Emily Denton. 2020. Saving face: Investigating the ethical concerns of facial recognition auditing. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. ACM, New York, NY, United States, 145–151.

[85] Eliane Röösli, Selen Bozkurt, and Tina Hernandez-Boussard. 2022. Peeking into a black box, the fairness and generalizability of a MIMIC-III benchmarking model. *Scientific Data* 9, 1 (2022), 1–13.

[86] Ognjen Rudovic, Jaeryoung Lee, Lea Mascarell-Maricic, Björn W Schuller, and Rosalind W Picard. 2017. Measuring engagement in robot-assisted autism therapy: a cross-cultural study. *Frontiers in Robotics and AI* 4 (2017), 36.

[87] Ognjen Rudovic, Yuria Utsumi, Jaeryoung Lee, Javier Hernandez, Eduardo Castelló Ferrer, Björn Schuller, and Rosalind W Picard. 2018. CultureNet: a deep learning approach for engagement intensity estimation from face images of children with autism. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, IEEE, New York, NY, United States, 339–346.

[88] Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet Vertesi. 2019. Fairness and abstraction in sociotechnical systems. In *Proceedings of the conference on fairness, accountability, and transparency*. ACM, New York, NY, United States, 59–68.

[89] Zhonghao Shi, Thomas R Groechel, Shomik Jain, Kourtney Chima, Ognjen Rudovic, and Maja J Matarić. 2022. Toward Personalized Affect-Aware Socially Assistive Robot Tutors for Long-Term Interventions with Children with Autism. *ACM Transactions on Human-Robot Interaction (THRI)* 11, 4 (2022), 1–28.

[90] Tom Simonite. 2015. Probing the dark side of google's ad-targeting system.

[91] Dimitris Spathis, Ignacio Perez-Pozuelo, Soren Brage, Nicholas J. Wareham, and Cecilia Mascolo. 2021. Self-Supervised Transfer Learning of Physiological Representations from Free-Living Wearable Data. In *Proceedings of the Conference on Health, Inference, and Learning* (Virtual Event, USA) *(CHIL '21)*. Association for Computing Machinery, New York, NY, USA, 69–78. https://doi.org/10.1145/3450439.3451863

[92] Statista. 2022. Global connected wearable devices 2016-2022 | Statista — statista.com. https://www.statista.com/statistics/487291/global-connected-wearable-devices/. [Accessed 14-Nov-2022].

[93] Angela Sucerquia, José David López, and Jesús Francisco Vargas-Bonilla. 2017. SisFall: A fall and movement dataset. *Sensors* 17, 1 (2017), 198.

[94] Harini Suresh and John Guttag. 2021. A framework for understanding sources of harm throughout the machine learning life cycle. In *Equity and access in algorithms, mechanisms, and optimization*. ACM, New York, NY, United States, 1–9.

[95] Sara Taylor, Natasha Jaques, Ehimwenma Nosakhare, Akane Sano, and Rosalind Picard. 2017. Personalized multitask learning for predicting tomorrow's mood, stress, and health. *IEEE Transactions on Affective Computing* 11, 2 (2017), 200–213.

[96] Oura Team. 2022. New to Oura: Blood Oxygen Sensing (SpO2) — ouraring.com. https://ouraring.com/blog/blood-oxygen-sensing-spo2/. [Accessed 14-Nov-2022].

[97] Kristina A Theis, Louise B Murphy, Dana Guglielmo, Michael A Boring, Catherine A Okoro, Lindsey M Duca, and Charles G Helmick. 2021. Prevalence of Arthritis and Arthritis-Attributable Activity Limitation—United States, 2016–2018. *Morbidity and Mortality Weekly Report* 70, 40 (2021), 1401.

[98] Connie W Tsao, Aaron W Aday, Zaid I Almarzooq, Alvaro Alonso, Andrea Z Beaton, Marcio S Bittencourt, Amelia K Boehme, Alfred E Buxton, April P Carson, Yvonne Commodore-Mensah, et al. 2022. Heart disease and stroke statistics—2022 update: a report from the American Heart Association. *Circulation* 145, 8 (2022), e153–e639.

[99] Yonatan Vaizman, Katherine Ellis, Gert Lanckriet, and Nadir Weibel. 2018. Extrasensory app: Data collection in-the-wild with rich user interface to self-report behavior. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. ACM, New York, NY, United States, 1–12.

[100] Michelle Van Ryn, Diana Burgess, Jennifer Malat, and Joan Griffin. 2006. Physicians' perceptions of patients' social and behavioral characteristics and race disparities in treatment recommendations for men with coronary artery disease. *American journal of public health* 96, 2 (2006), 351–357.

[101] Dimitrios Vasdekis, Sofia Yfantidou, Stefanos Efstathiou, and Athena Vakali. 2022. WeMoD: A Machine Learning Approach for Wearable and Mobile Physical Activity Prediction. In *2022 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops)*. IEEE, IEEE, New York, NY, United States, 385–390.

[102] Praveen Veerabhadrappa, Matthew Duffy Moran, Mitchell D Renninger, Matthew B Rhudy, Scott B Dreisbach, and Kristin M Gift. 2018. Tracking steps on apple watch at different walking speeds. *Journal of general internal medicine* 33, 6 (2018), 795–796.

[103] Sahil Verma and Julia Rubin. 2018. Fairness definitions explained. In *2018 ieee/acm international workshop on software fairness (fairware)*. IEEE, IEEE, New York, NY, United States, 1–7.

[104] Guang Wang, Yongfeng Zhang, Zhihan Fang, Shuai Wang, Fan Zhang, and Desheng Zhang. 2020. FairCharge: A data-driven fairness-aware charging recommendation system for large-scale electric taxi fleets. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 1 (2020), 1–25.

[105] Rui Wang, Fanglin Chen, Zhenyu Chen, Tianxing Li, Gabriella Harari, Stefanie Tignor, Xia Zhou, Dror Ben-Zeev, and Andrew T Campbell. 2014. StudentLife: assessing mental health, academic performance and behavioral trends of college students using smartphones. In *Proceedings of the 2014 ACM international joint conference on pervasive and ubiquitous computing*. ACM, New York, NY, United States, 3–14.

[106] Yifan Wang, Weizhi Ma, Min Zhang*, Yiqun Liu, and Shaoping Ma. 2022. A survey on the fairness of recommender systems. *ACM Journal of the ACM (JACM)* 111 (2022), 1–43.

[107] Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. 2020. Towards fairness in visual recognition: Effective strategies for bias mitigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. IEEE, New York, NY, United States, 8919–8928.

[108] Jenna Wiens, W Nicholson Price, and Michael W Sjoding. 2020. Diagnosing bias in data-driven algorithms for healthcare. *Nature medicine* 26, 1 (2020), 25–26.

[109] Samuel Yeom and Michael Carl Tschantz. 2021. Avoiding disparity amplification under different worldviews. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. ACM, New York, NY, United States, 273–283.

[110] Sofia Yfantidou, Pavlos Sermpezis, and Athena Vakali. 2021. Self-tracking technology for mhealth: A systematic review and the past self framework. *arXiv preprint arXiv:2104.11483* (2021), 1–40.

[111] Guo-Qiang Zhang, Licong Cui, Remo Mueller, Shiqiang Tao, Matthew Kim, Michael Rueschman, Sara Mariani, Daniel Mobley, and Susan Redline. 2018. The National Sleep Research Resource: towards a sleep data commons. *Journal of the American Medical Informatics Association* 25, 10 (2018), 1351–1358.

[112] Haoran Zhang, Amy X Lu, Mohamed Abdalla, Matthew McDermott, and Marzyeh Ghassemi. 2020. Hurtful words: quantifying biases in clinical contextual word embeddings. In *proceedings of the ACM Conference on Health, Inference, and Learning*. ACM, New York, NY, United States, 110–120.

## A    FAIRNESS TAXONOMY IN MACHINE LEARNING

Viewed through the lens of quantitive science, machine learning research has broadly grouped fairness into two categories: *individual fairness* and *group fairness*. In the broad sense, group fairness partitions the general population into groups based on sensitive (a.k.a. protected) attributes and seeks statistical equality across groups. On the other hand, individual fairness seeks for similar individuals to be treated similarly [35, 75].

In individual fairness, determining whether individuals are similar requires first defining what features are relevant to fairness [41]. However, in PI, it would be incomplete to define such similarity solely based on digital behavioral biomarkers, such as steps, or heart rate, as hidden contextual information might be significantly more relevant. To this end, given access to de-identified aggregated information, we proceed with group fairness metrics and definitions from now onward. This decision translates to our use case as exploring significant differences in allocation, representation, or error rates regarding future step goals across different population segments. For example, do females get systematically lower step goals than their male counterparts?
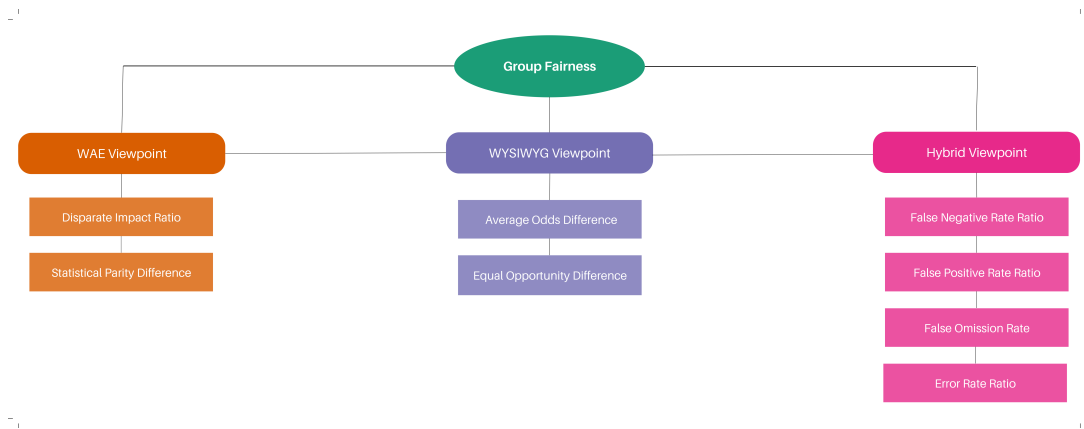


Fig. 12. An organization of fairness metrics in machine learning according to different worldviews: WAE, WYSIWYG, and a hybrid between the two.

Within group fairness, there are still two opposing viewpoints: we're all equal (WAE) and what you see is what you get (WYSIWYG) [44, 109]. The WAE viewpoint considers that all groups have similar abilities to perform the task, e.g., all groups of people are equally capable of walking more, while the WYSIWYG viewpoint holds that the data reflect each group's ability to perform the task, e.g., some groups of people might be less capable of walking more. Group fairness metrics lie under either viewpoint or somewhere in between, as seen in Figure 12. Overall, every metric tries to quantify -one way or another- the difference in performance between privileged and unprivileged groups of users. We provide detailed definitions, mathematical formulas, PI-specific interpretations, and visual representations of multiple fairness metrics in Appendix B. All metrics' definitions and formulas are taken from the AI Fairness 360 (AIF360) toolkit [20].

## B    GROUP FAIRNESS METRICS & INTERPRETATION

This section presents the most common fairness metrics across all viewpoints, namely WAE, WYSIWYG, and hybrid. For each metric, we provide a short definition, a mathematical formulation, and the metric's usage and bias interpretation concerning our use case. The confusion matrix (see Table 3) is the heart of performance measurement in machine learning and is also used in the fairness metrics definitions below. All metrics are

expressed as *ratios* or *differences* between unprivileged (u) and privileged (p) groups. Note that $D$ is the user sample, $\hat{Y}$ is the predicted label, and pos_label is what we consider as a positive outcome scenario (e.g., high physical activity). Below, we present the metrics per viewpoint.

Table 3. The confusion matrix for performance measurement of models, including standard metrics.

| | | **True** | | |
| --- | --- | --- | --- | --- |
| | | Positive Label | Negative Label | |
| **Predicted** | Positive Label | TP | FP | False Discovery Rate<br>FDR = $\frac{FP}{FP+TP}$ |
| | Negative Label | FN | TN | False Omission Rate<br>FOR = $\frac{FN}{FN+TN}$ |
| | | False Negative Rate<br>(a.k.a Miss-rate)<br>FNR = $\frac{FN}{TP+FN}$ | True Negative Rate<br>(a.k.a. Specificity)<br>TNR=$\frac{TN}{FP+TN}$ | |
| | | True Positive Rate<br>(a.k.a. Sensitivity)<br>TPR = $\frac{TP}{TP+FN}$ | False Positive Rate<br>(a.k.a. Fall-out)<br>FPR = $\frac{FP}{FP+TN}$ | |

## B.1   WAE Metrics

The WAE viewpoint supports that the data, e.g., measured physical activity, may contain biases -this holds true for physical activity, as we will discuss in a later section-, so their distribution being different across groups should not be mistaken for a difference in ability. The two most commonly used WAE metrics are the *Disparate Impact Ratio (DIR)* and the *Statistical Parity Difference (SPD)* (See Table 4).

## B.2   Hybrid Metrics

Hybrid metrics lie in-between the two viewpoints. The most commonly used hybrid metrics, depending on the problem's context, are the *False Positive Rate (FPR) Ratio*, the *False Negative Rate (FNR) Ratio*, the *False Omission Rate (FOR) Ratio*, and the *Error Rate Ratio (ERR)* (See Table 5).

Table 4. WAE Metrics' definitions, formulas, and task and bias interpretations specific to our use case.

| Metric | Definition | Formula | Task Interpretation | Bias Interpretation |
|---|---|---|---|---|
| DIR | The ratio of base or selection rates between unprivileged and privileged groups. | $\dfrac{\Pr(\hat{Y} = \text{pos\_label} \mid D = u)}{\Pr(\hat{Y} = \text{pos\_label} \mid D = p)}$ | How many users receive high activity goals in the unprivileged group compared to the privileged group? | A low *DIR* (*DIR* < 1) indicates that the unprivileged user group systematically receives fewer high activity goals. |
| SPD | The difference in selection rates between unprivileged and privileged groups. | $\Pr(\hat{Y} = \text{pos\_label} \mid D = u) - \Pr(\hat{Y} = \text{pos\_label} \mid D = p)$ | Same as above | A low *SPD* (*DIR* < 0) indicates that the unprivileged user group systematically receives fewer high activity goals. |

## B.3 WYSIWYG Metrics

The WYSIWYG viewpoint supports that the data, e.g., measured physical activity, correlates well with future activity and that there is a way to use them correctly to compare the abilities of the users. The two most commonly used WAE metrics are the *Average Odds Difference (AOD)* and the *Equal Opportunity Difference (EOD)*.

## B.4 Fairness Metrics' Value Ranges

Figure 13 gives a graphical overview of all fairness metrics discussed. We notice that difference-based metrics have a different range of values compared to ratio-based metrics. Specifically, difference-based metrics are within the [−100%, +100%] range, while ratio-based metrics are within the [0, ∞] range. For both categories, a value of 1.0 is optimal, indicating demographic parity. Anything greater or less than the optimal value indicates some level of bias. According to AIF360, accepted difference-based metrics' values are within the range [−0.1, +0, 1], and accepted ratio-based metrics' values are within [0.8,1.25], but such ranges are not universally accepted and might be adjusted on a task-by-task basis. Table 7a indicates which user group, namely privileged or unprivileged, benefits based on a group fairness metric's value.

Table 5. Hybrid Metrics' definitions, formulas, and task and bias interpretations specific to our use case.

| Metric | Definition | Formula | Task Interpretation | Bias Interpretation |
|---|---|---|---|---|
| FPR Ratio | The ratio between the number of negative outcomes wrongly categorized as positive, i.e., false positives (FP), and the total number of actual negative outcomes regardless of classification. | $\dfrac{FPR_{D=u}}{FPR_{D=p}}$ | From all the low active users, how many wrongfully received high activity goals? | A low *FPR* Ratio (FPR Ratio < 1) indicates that the privileged low active user group systematically receives more high activity goals compared to the unprivileged low active user group. |
| FNR Ratio | The ratio between the number of positive outcomes wrongly categorized as negative, i.e., false negatives (FN), and the total number of actual positive outcomes regardless of classification. | $\dfrac{FNR_{D=u}}{FNR_{D=p}}$ | From all the highly active users, how many wrongfully received low activity goals? | A high *FNR* Ratio (FNR Ratio > 1) indicates that the unprivileged highly active user group systematically receives more low activity goals compared to the privileged highly active user group. |
| FOR Ratio | The ratio between the outcomes wrongly categorized as negative, i.e., FN, and the total number of classified negative outcomes. | $\dfrac{FOR_{D=u}}{FOR_{D=p}}$ | From all the users that were given low activity goals -rightfully so or not-, how many were actually highly active? | A high *FOR* Ratio (FOR Ratio > 1) indicates that the unprivileged user group systematically receives more wrong low activity goals compared to the privileged user group. |
| ERR | The ratio between the erroneous outcomes, i.e., FN or FP, and the total number of outcomes. | $\dfrac{ER_{D=u}}{ER_{D=p}}$, where $ER = \dfrac{FP + FN}{P + N}$ | How many times was the activity level prediction wrong? | A high *ERR* (ERR > 1) indicates that the unprivileged user group systematically receives more wrong goals -low or high- compared to the privileged user group. |

Table 6. WYSIWYG Metrics' definitions, formulas, and task and bias interpretations specific to our use case.

| Metric | Definition | Formula | Task Interpretation | Bias Interpretation |
|--------|-----------|---------|---------------------|---------------------|
| EOD | The difference of true positive rates between the unprivileged and the privileged groups. | $TPR_{D=u} - TPR_{D=p}$ | From all the highly active users, how many were actually given high activity goals? | A low *EOD* ($EOD < -0.1$) indicates that the unprivileged highly active user group systematically receives fewer high activity goals compared to the privileged highly active user group. |
| AOD | The average difference between the FPR and the TPR between unprivileged and privileged groups. | $\frac{(FPR_{D=u} - FPR_{D=p}) + (TPR_{D=u} - TPR_{D=p})}{2}$ | TBD | TBD |

Table 7. Value interpretation for group fairness metrics. The table shows which user group is treated unfairly -in a negative manner- in each case. UN indicates the unprivileged user group, and PR indicates the privileged user group. UN ~ PR indicates a fair outcome. Notice that the same values may mean different things in the case of ratio-based metrics.

(a) Ratio-based metrics.

| Metric | $v < 0.8$ | $0.8 < v < 1.25$ | $v > 1.25$ |
|--------|-----------|------------------|------------|
| *DIR* | UN | UN ~ PR | PR |
| *FOR* | PR | UN ~ PR | UN |
| *FNR* | PR | UN ~ PR | UN |
| *FPR* | UN | UN ~ PR | PR |
| *ERR* | PR | UN ~ PR | UN |

(b) Difference-based metrics.

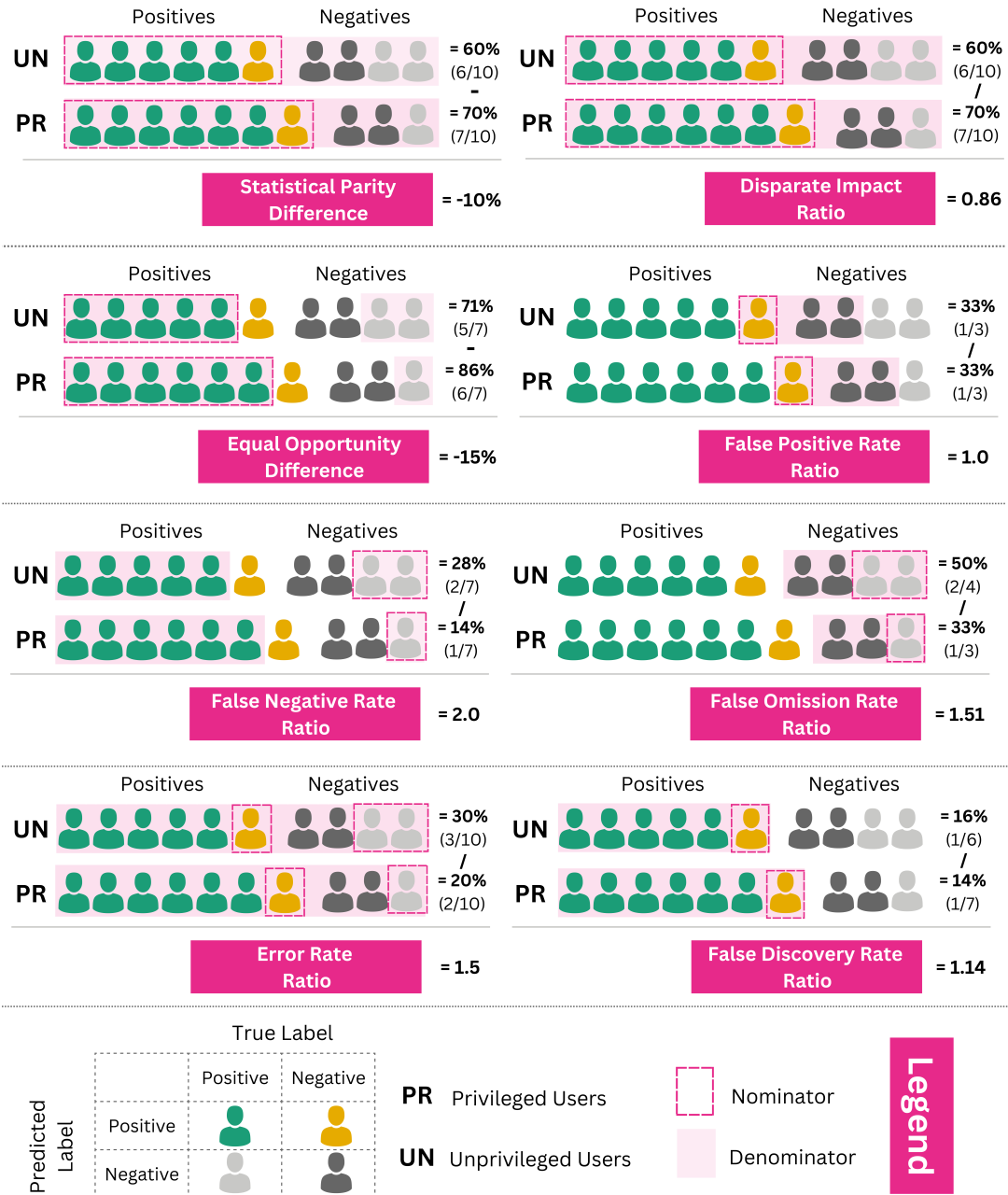| Metric | $v < -0.1$ | $-0.1 < v < 0.1$ | $v > 0.1$ |
|--------|-----------|------------------|-----------|
| *SPD* | UN | UN ~ PR | PR |
| *EOD* | UN | UN ~ PR | PR |
| *AOD* | UN | UN ~ PR | PR |

Fig. 13. A graphical overview of the fairness metrics discussed.