

# Ethical trade-offs in medical Machine Learning technologies

Albrecht, Thomas (5733587)      Petruck, Julian (3857386)  
Jaques, Arthur (5998179)

March 6, 2022

## Abstract

We use the application of Machine Learning (ML) to healthcare as a case study of ethical trade-offs. We first introduce the concept of Pareto optimality to formalize trade-offs in general. Then we identify three main trade-offs in fair ML: The choice of fairness measure and the trade-offs between accuracy, fairness and privacy. We also discuss some ideas that have been introduced to handle them. We examine the particular trade-offs that appear in medical practice, be it in policy, treatment, or research. Because there is psychological resistance to trade-off talk in health care since it touches on the ‘sacred value’ of life, such trade-offs are often identified as moral dilemmas. We show that fairness problems have been identified in unaided medicine, and many issues identified in the fair ML literature are just their derivations. After arguing for the need of ethical assessments of ML tools that are relative to clinicians, we analyze what differences ML decision support tools have to human doctors with respect to fairness, and what they can positively contribute. Finally, we look beyond the issue of fairness towards other ethical trade-offs. Our analysis shows that while very often trade-offs at the intersection of the fields of medicine and ML look like new problems, most of them can be related to trade-offs we are used to dealing with in other technologies.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Introduction to trade-offs</b>	<b>2</b>
2.1	Multi-objective optimization and Pareto optimality . . . . .	2
2.2	Trade-offs in fair Machine Learning . . . . .	4
<b>3</b>	<b>Trade-offs in medicine</b>	<b>6</b>
3.1	Trade-offs in traditional medicine . . . . .	6
3.2	Fairness in traditional medicine . . . . .	7
3.3	Solutions from medical ethics . . . . .	7
<b>4</b>	<b>Fairness in medical Machine Learning systems</b>	<b>8</b>
4.1	Fairness problems . . . . .	8
4.2	Domain-dependence of fairness problems . . . . .	9
4.3	Potential benefits of Machine Learning, and the sin of perfection . . . . .	10
<b>5</b>	<b>Further ethical trade-offs in medical Machine Learning systems</b>	<b>13</b>
5.1	Trade-offs . . . . .	13
5.2	We have seen that before . . . . .	15
<b>6</b>	<b>Conclusions</b>	<b>17</b>

# 1 Introduction

Machine Learning (ML) in health care and medicine has grown to be one of the most discussed, but also most promising applications of the ever-growing technology of ML. In recent years, more and more research has shown ML to be an effective way of supporting health care practitioners in a great diversity of ways [39, 46]. However, there is also growing concern about the implications the deployment of ML has for the future of health care and medicine.

While there are ongoing public discussions about ML replacing humans as workers in different ways [40], many scholars have made clear that ML and AI tools will not replace clinicians in the near future, but rather be integrated as support systems, for example as clinical decision support systems (CDSS) [36].

CDSS have been used since the 1980s with growing success [44]. Only in recent years, the involvement of Machine Learning in those systems has led to a new regulatory situation. Still, their deployment and success can tell us a lot about the way to go with ML tools. For example, although closed loop systems, i.e., systems where every step of the process from diagnosis to drug intake is computerized and monitored, do already exist, they are not commonly used, partly due to costs but certainly also due to the involved surveillance environment for patients [44].

We consider CDSS a particularly interesting use case of ML when it comes to fairness and more generally ethical discussion, since the biomedical field historically played an important role when it comes to ethical principles and deliberations [19, 47]. It helped displace purely theoretical ethical discussions (meta-ethics) to more concrete, unavoidable, and tangible questions (applied ethics). It further helped surpass purely relativist, subjectivist, and psychological investigations of ethics [47]. Broad, universal moral principles were replaced by case studies, for example arising in clinical medicine. Finally, ethical considerations started taking into account the roles and relationships of the actors present (for example, by recognizing the authority relationship between doctors and patients). Through all those changes that it instigated, poetically put, “[m]edicine saved the life of ethics” [47].

A further reason of interest is that medicine poses some inevitable moral dilemmas that can not be ignored, about which a substantial ethical literature has been developed. The high stakes found in some decisions in medicine might induce less readiness to sacrifice performance for fairness. Additionally, in health care settings we are often interested in the prevention of harm rather than the allocation of goods (such as in job placements, college admissions, and other areas where fair-ML is being applied). This means that solutions such as the randomization of a group’s predictions to ensure equal harms are considered unacceptable. A final point of interest is the fact that in medicine, sensitive attributes such as gender and race might be true predictors. Hence, avoiding group differences would harm everyone [34], which makes solutions such as ‘fairness through unawareness’ and independence-based approaches very questionable.

In our investigations, we draw on literature from the fields of fair-ML, economics, medicine, biomedical ethics, and philosophy. In Chapter 2, we examine economists’ work

on the notion of trade-offs to derive a meaningful formalization of them. We then examine trade-offs that are frequently discussed in the ML literature. The analysis of the medical literature found in Chapter 3 helps us identify the trade-offs (often described as moral dilemmas) inherent to the field of health care. We furthermore examine fairness problems in medical practice, and discuss the usefulness of principled approaches found in the bioethical literature in enlightening the ethical discussion about CDSS. In Chapter 4, we relate fairness problems raised by the fair-ML community to ethical questions raised in the medical field. In particular, we argue that many fairness problems are inherent to medical practice and not introduced by the application of ML. We then examine what ML can positively contribute with respect to ethical issues in medicine. Finally, in Chapter 5 we show how the application of ML to health care introduces new trade-offs and ethical questions. We analyze the supporting literature critically, by discussing in which measure those questions can be found in ethical discussions about other technologies.

## 2 Introduction to trade-offs

The notion of a trade-off describes a decision between multiple (usually mutually contradictory) objectives, in the sense that a gain in one objective results in loss in one or more other objectives. On a broad view, trade-offs are the basic problem of human governance (“the central rationale for many policies” [19, p. 77]): How many resources we allocate for one problem, leaving less for another one. Trade-offs are intuitively understood from a young age, as they are very common in everyday life, and encompass all human decision-making. Biology, evolutionary theory, and the human body can all be understood in terms of trade-offs [30]. Similarly, policies dealing with large-scale stochastic problems (vaccination, traffic security, nuclear deterrence, criminal justice) always entail harms and benefits [19], and hence trade-offs.

But in economics in particular trade-offs are of special interest, making them a central focus of study in the field. Accordingly, economists have proposed multiple approaches to formalize them. One such approach, which is so widespread and commonly used that it can be regarded as a convention, is Pareto efficiency and the Pareto front.

### 2.1 Multi-objective optimization and Pareto optimality

To approach the choice of an optimal feasible decision (allocation) for various types of trade-offs we introduce multi-objective optimization. A general multi-objective optimization problem  $F$  can be written as a maximization in the following way:

$$\max F(x) = (u_1(x), \dots, u_k(x)), \quad s.t. \ x \in X$$

Here  $X$  denotes the set of all feasible decisions and  $u_i(x)$  are the utility/objective functions, adding up to  $k$  dimensions. For a non-trivial multi-objective optimization problem it is not possible to maximize every single objective function at the same time. Thus the notion

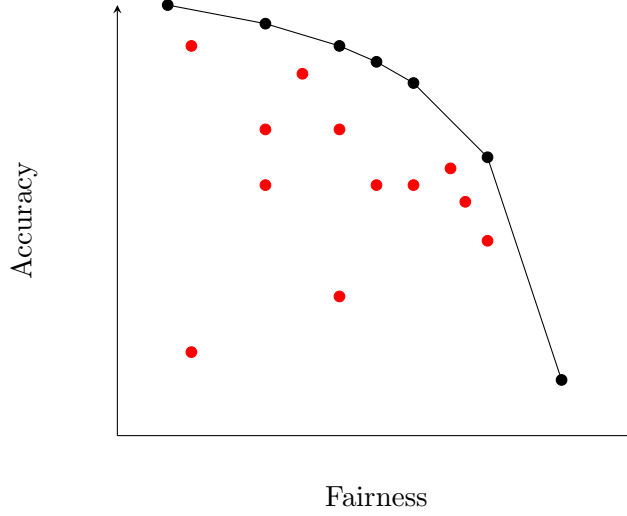


Figure 1: Fictional example of a two-dimensional Pareto front (black). Each point is a solution in the feasible set, the red points are dominated. Along the front increasing one objective leads to decreasing the other one.

of Pareto optimality is introduced: A decision  $x \in X$  is said to Pareto dominate another solution  $x' \in X$  if the following both hold:

1.  $\forall i \in 1, 2, \dots, k : u_i(x) \geq u_i(x')$
2.  $\exists j \in 1, 2, \dots, k : u_j(x) > u_j(x')$

A non-dominated decision is also called Pareto optimal or Pareto efficient. Any Pareto optimal decision cannot be further improved for one objective unilaterally without resulting in loss in one or more other objectives. The set of all Pareto optimal decisions is called the Pareto front. Note that Pareto optimality doesn't ensure anything beyond the property derived above. In particular it doesn't provide any guarantees about a fair or normative allocation or decision [49]. If the optimization problem is two-dimensional the Pareto front can be visualized in an intuitive way (see Figure 1): The objectives are the axes in a 2D plane, moving along the Pareto front showcases intuitively how increasing one objective decreases the other one.

Of course, it would be desirable to avoid the many trade-offs, in the sense of maximizing all objectives simultaneously. Unfortunately, this maximizing solution might not be in the feasible set. But the concept of Pareto optimality alone won't result in a single optimal or 'best' answer to our decision problem. Rather, the approach eliminates all 'strictly worse' possible decisions in the feasible set and the decision maker is faced with a new problem, they now have to choose one solution (decision) from the Pareto front. Luckily, the search space is (usually) much more narrow when only considering the Pareto optimal solutions.

Also, the trade-offs are made explicit, in the sense that choosing to improve one objective along the frontier will incur ‘cost’ in one or multiple other objectives. Usually the decision maker needs to choose a single decision from the Pareto front. Depending on the problem at hand the decision maker could (or rather has to) potentially incorporate additional prior information (expert knowledge or preference). This information can be incorporated a priori into the global utility function (scalarization), used interactively during calculation of the front, or applied after the front has been calculated to choose a solution (a posteriori approach) [24].

## 2.2 Trade-offs in fair Machine Learning

When designing any technology [1] there are many trade-offs inherent in the process. Of course ML systems and algorithms are no exception. For example, in medical ML systems we aim to guarantee both privacy and fair treatment, while our primary goal is still to achieve high accuracy and avoid misclassifications or mistreatment [7]. This leads to multiple trade-offs that have to be made. Moving forward we are going to characterize three main fairness related trade-offs we identified in machine learning systems:

**Choice of a fairness measure** Often we will quantify fairness using a selected fairness measure. Some often used measures are equalized odds, statistical parity and predictive parity [18]. The factors to take into account when deciding on what metric of fairness to use for a specific ML system are multiple. The different metrics each formalize different notions of morality [5]. First we need to decide what moral principles we want to follow, i.e. , what we mean by equal or just treatment. But the choice of measure entails a trade-off already [3, 8, 29], as some notions of fairness are mutually contradictory and cannot be satisfied at the same time. As a result, we have to choose between several contradictory measures, given that they are feasible for the ML problem at hand. All the metrics mentioned above are group fairness measures. In group fairness we try to protect so-called sensitive attributes, in the sense of treating different groups (distinguished by the values for the sensitive attributes) equally. But how do we choose the sensitive features or groups we want to protect? We cannot take the sensitive features for granted, instead the separation into groups is also always a choice that humans have to make, based on legal and societal considerations among others.

One apparent solution to the problem of choosing a group fairness measure could be to instead consider so-called individual fairness. Here we consider equal treatment of similar individuals, independent of any group affiliation. By only considering individuals instead of groups we seemingly sidestep the problem of choosing the sensitive features we want to protect and which fairness measure to use. Choosing an appropriate measure for the similarity of individuals poses difficulties, as there are many feasible measures, and we cannot restrict the space of metrics by imposing additional (moral) constraints on the choice of the fairness measure. But certain similarity measures do correspond to certain moral notions of equal treatment, e.g. statistical parity [15].

This leads to a conflict between individual fairness, with individuals wishing to be judged independently of their group identity, and group fairness, which tries to correct for discrimination based on sensitive features.

**Accuracy and Fairness** Prediction accuracy is a very desirable property in ML systems, maximizing it often being the primary goal of the employed algorithm. Fair machine learning is concerned with identifying and mitigating bias and discrimination of sensitive attributes in ML systems. Ideally we would like to achieve optimal accuracy while not discriminating with respect to any sensitive feature. But as demonstrated empirically in e.g. [26] and [54] avoiding discrimination (or achieving a certain level of fairness) often directly results in the loss of prediction accuracy.

Furthermore, [35] showcases that this trade-off is a property inherent in the data and cannot be avoided by a clever choice of the learning algorithm used when learning on a modified problem subject to a fairness constraint. They show that the accuracy and fairness trade-off depends on the ‘alignment’ of the label and the sensitive feature, in the sense that if the label and sensitive feature are highly correlated, ensuring a certain level of fairness will result in huge loss of accuracy. Conversely, if the sensitive feature and the label are perfectly independent of each other, we can achieve perfect fairness while retaining the full accuracy.

To guarantee a certain amount of fairness, subject to a certain fairness measure, while retaining the maximum accuracy possible under that fairness constraint we can examine the Pareto front of the trade-off. This is characterized by [31] and [51], where the trade-off between accuracy and fairness is given as a Pareto front for different measures of fairness. This allows for examination and comparison of the nature of the trade-off for different problems and measures.

**Accuracy and Privacy** Privacy, just like fairness, is another information-based harm, albeit with slightly different characteristics [48].

In the simplest approach to the problem, to assure a certain degree of privacy we have to discard or randomize some of the data, which in turn will lead to worse accuracy. Multiple more refined approaches have been proposed to mitigate this loss of accuracy (many approaches for privacy preserving data mining and ML e.g. [14], which uses differential privacy).

But while those are able to preserve accuracy to some degree under certain privacy guarantees, they aren’t always applicable and don’t necessarily resolve all of the harms mentioned. This can be partially explained by the different meanings of ‘privacy’ in different fields. Differential privacy might additionally amplify the unfairness of a given model [2], or in some cases even be incompatible with fairness in general [11]. This is a problem, as it is often desirable to achieve high accuracy while preserving both privacy and fairness, e.g. when we try to balance all three axes of conflict (accuracy, privacy and fairness) in our ML system [7].

### 3 Trade-offs in medicine

Let us now consider trade-offs that can be observed in the practice of medicine. In this chapter, we concentrate on trade-offs found before the application of Machine Learning.

#### 3.1 Trade-offs in traditional medicine

As hinted earlier, trade-offs play a central role in medicine. They can appear at different levels, be it in high level healthcare policy decisions or when considering different possible treatment options for one specific patient. Mentioning trade-offs in the field of medicine might cause defensive reactions because of the mathematical flavor the notion carries, which clashes with the supposed complexity of ethical problems [52]. Furthermore, the simple concept of trade-offs in cases where sacred values (such as human lives) clash with secular values (such as money) is typically morally disturbing to the public, and thus avoided [45]. Suggesting that doctors apply moral trade-offs in their practice is a contestable affirmation, since the nature of their ethical deliberations is necessarily partly non-mathematical [53]. Hence, a more accepted term here is ethical (or moral) dilemma, defined as a problem that arises when opposing values or principles co-occur [42]. Fundamentally, however, trade-offs and practical solutions to moral dilemmas are the same thing: Decisions on how much to weigh principles that can not be fully respected at the same time. If for example the administration of a treatment could harm the patient as a side effect, one might still choose to treat the disease if it is the lesser evil (as is the case with chemotherapy [37]). In this situation, the physician faces multiple trade-offs: He has to consider the effectiveness of the administered treatment (which is uncertain for the given patient), the likelihoods, and the magnitude of possible adverse side effects (which are also uncertain for the given patient). Far from purely qualitative reasoning, a step in the quantitative dimension of trade-offs is shown for example by evidence-based medicine [30], which serves to inform decisions on what risks are to be taken with the promise of some potential benefit.

The perhaps most obvious trade-off in the practice of medicine, that every doctor understands, is the one between the potential gains and the risked losses of a given treatment [30]. In fact, one can go as far as to “conceptualize medicine itself as the art of managing trade-offs” [30, p. 575]. From the doctor’s allocation of time to specific patients, to the risk of switching to a new potentially better treatment, to the decision of how aggressively to treat terminal patients, every hard decision a medical practitioner has to take entails a trade-off.

Health care itself, as part of governance, is ridden with trade-offs. Health care systems are administrated according to risk-benefit analyses, both as part of the overall governance budget (how many resources to invest in health care as opposed to, for example, transportation) and within the system (which operations to prioritize and what costs to cover among other questions) [12]. Empirical research suggests that the phrasing of such decisions has a big impact in the public’s perception of the problem: “Hospital administrators wrestling with tragic trade-offs can find themselves in the dock as soon as critics wonder who set the budget constraint that made it possible to save only one child” [45, p. 323]. This might



explain the resistance to trade-offs talks in health and preference for moral dilemma formulations mentioned above, despite the arguably higher practical utility of the concept of trade-off.

Furthermore, much attention is paid to preserving privacy when using medical records and clinical data for scientific studies. The European GDPR is for example an important personal data protection law, that because of unclarity and unresolved legal issues often stalls scientific research and progress [16]. A balance between the protection of the personal sensitive data of patients and potential scientific advancements must thus be found, leading to a necessary trade-off.

A more hidden trade-off, masked by naive claims of the complete objectivity of science, is the one between methodological criteria in clinical research. With the modern focus on evidence-based medicine and randomized controlled trials as research instruments, and the consequent diminution in the perceived value of cohort studies, case-control studies, expert opinion, and case studies, active (interventional) studies have been elevated to the golden standard of medical research. However, in doing so, the whole focus is placed on the methodological criteria of generality and precision, disregarding criteria such as realism, coherence, explanatory power, and others. This hides the underlying trade-off between methodological criteria, giving an absolute choice where case-by-case considerations of medical focus and contextual values are required to set methodological priorities [22].

As mentioned before, ML often has to deal with trade-offs between fairness and performance. Let us in the following examine what fairness problems can arise in medical practice, that will need to necessary trade-offs.

### **3.2 Fairness in traditional medicine**

We can identify some key fairness issues in the medical literature. Healthcare disparities are a well-accepted reality, and “often encompass all 5 domains of the social determinants of health as defined by the US Department of Health and Human Services (economic stability, education access and quality, healthcare access and quality, neighborhood and built environment, and social and community context)” [6, p. 2]. For example, gender bias is a recognized factor in health care. It is observed e.g. in the uneven composition of clinical trials samples, concentration on male-typical risk factors in studies, and in the different seriousness with which men’s and women’s complaints are received by medical doctors [41]. This bias might pose fairness problems, since the resulting health care system might overperform on men and underperform on women.

### **3.3 Solutions from medical ethics**

A good starting point for ethical discussions in medicine are the well-established guiding principles in biomedical ethics proposed by Beauchamp and Childress: Respect for autonomy, beneficence, non-maleficence and justice [36, 38, 42]. The guiding ideas of biomedical ethics can be used to assess specific applications of ML to health, for example decision

support in occupational health, by considering the potential benefits and risks with respect to those principles [42].

Despite the interest of considering biomedical ethics, actual practice seems to indicate that case-by-case evaluations of the moral implications of medical decisions are more useful than principled approaches. Toulmin[47] reports for example how a commission of people from different backgrounds, faced with specific practical problems, were able to reach some consensus (disagreeing at most about the degree of the decisions), all while furiously disagreeing about the principles supporting their decisions [47]. It must furthermore be noted that physicians typically exert their clinical judgment only after collecting a precise case history, instead of following general theoretical considerations early on [47]. So if we take their behavior as somewhat exemplary, case-by-case evaluations considering the specific details might be required.

## 4 Fairness in medical Machine Learning systems

Having identified a number of fairness issues entailed in the practice of medicine, we now turn to problems that were identified by the fair-ML community when developing support systems. We first provide a short summary, and then argue that most of the emerging concerns are not fundamentally new, but related to older ethical problems intrinsic to (medical) decision problems. After arguing that the performance and fairness of CDSS should be analyzed not absolutely, but relatively to what human practitioners currently do, we evaluate what contributions and improvements (as it relates to fairness) ML could bring to medicine.

### 4.1 Fairness problems

Ethical questions about algorithmic unfairness are a type of normative concerns [36]. The general concern is that CDSS trained on unbalanced or biased datasets might pick up the wrong patterns and exacerbate existing inequalities in health by overfitting on advantaged groups [6, 36]. The problem is often identified in the data used for training, which might contain label prejudice (a kind of negative legacy), variability amongst clinicians and institutions, and evolving clinical knowledge [6]. The data at deployment time is also identified as a source of problems: Population shift makes the developed model not adapted to the current population, and the usual lack of ground truth labels at test time makes the evaluation difficult [6]. Ethical discussion emerges in classification problems, when the operating point of the algorithm is chosen [17]. This corresponds to the trade-off between false positives and false negatives, a well-known question in medicine (for example in breast cancer screening [17]). Other examples of biases of interest are sampling bias, unforeseen correlations, true systemic bias with biological causes, and batch effects [17], as well as negative legacy and labeling prejudice [6]. Furthermore, it must be recognized that “algorithmic development is never an entirely objective, value-free endeavour: it will be influenced by a host of social and institutional norms, practices and attitudes that could well build bias

into design.” [53, p. 673]. While bias is an ethically neutral term indicating statistical imbalances, unfairness is the judgement of bias as incompatible with moral principles [17]. Fairness in ML is usually defined in terms of groups, quantified by group fairness measures (refer to 2.2 for details). The consensus is that fairness through unawareness (removing group information from the training features) is not the right solution, because of sensible attribute leakage and the true effects of latent biological factors in many diseases [6].

## 4.2 Domain-dependence of fairness problems

The practical problem of applying ML to health care decision tasks carries with it a certain number of unavoidable decisions about the relative importance of contrasting principles. Furthermore, imposing fairness constraints typically leads to lower performance of the ML algorithms. For these reasons, trade-offs have to be made (refer to 2.2). Above, we concentrated on trade-offs related with fairness considerations. An interesting aspect of such trade-offs is their origin. We argue that many fairness-related problems, which lead to trade-offs in their solutions, originate from the medical (decision) problem itself, and not from the technology used to solve it. This means that the analysis of fairness is necessarily domain-dependent [17], and in our case should take inspiration from medical ethics literature [36]. Additionally, “considerations internal to medical science” and “contextual values” must help inform the decision on which methodological criteria to focus [22, p. 252], and whether the available data and ML technology allow such a focus. Hence, we suggest abstracting from ML itself and regarding it as a general technology used to solve a pre-existing problem. We take inspiration from philosophy of technology, recognizing the interaction between our behavior and the technology we use [42]. Our reason to, on a first analysis, disregard the specifics of ML is that in many cases, ML often simply specifies existing trade-offs and makes them unavoidable.

As a concrete example, consider the problem of unbalanced training data, causing the ML algorithm to reach better performance on over-represented groups. This problem is hidden, but still present, in unassisted medicine. Medical practitioners arguably learn the best treatments partly from experience [53]. If the majority of their patients come from a particular group, it is very likely that they will ‘overfit’ their knowledge to that group, or at least be able to predict their response to particular treatments better. Similarly, textbook knowledge is partly derived from observations of medical practitioners and/or statistical studies. Biased data informing those studies will bias the observed results [6]. This effect is exemplified by heart health research, where research on biased data (higher proportion of men) lead to uneven improvements in heart attacks treatment with respect to gender [34]. The data informing that knowledge is probably very similar to the data used to train ML algorithms. The resulting inferences will hence be similarly biased as a consequence of the shown data imbalances. In this regard, ML systems have a better potential to solve the problem, for example by using importance weighting or under-/over-sampling [6], or by driving the collection of more diversified data [53]. By no means do we intend to suggest that the solution is easy, since blindly applied technical fixes may introduce undetected harms (contrasting with the bioethical principle of non-maleficence). However, ML has the

potential to fix biases in a way that traditional medical practice can not [34] (although traditional statistics can help, see e.g. importance sampling).

Technical bias, a consequence of new knowledge or data not being integrated into the algorithmic system, is arguably similar to emergent bias that is typical of any decision system. Human deciders are for example subject to the ‘availability heuristic’, and are often mandated continuing education to keep their knowledge up to date, to avoid emergent bias [53]. Similar mandated technical solutions (software upgrades, maintenance procedures) are in principle possible fixes for the problem when ML systems are used [53].

The issue of reducing an individual to a group identity already exists in statistics, and arises in classical clinical practice as well. As an example for why this is considered a problem, there is evidence for the strong moral resistance to the use of statistics (such as generalization-allowing base rates) in sensitive situations [45]. Despite claims about the possibility of ‘personalized medicine’ thanks to ML methods, individuals are still reduced to their features [4, 42]. Although it is true that personalization might be an illusion, ML allows for more precise groupings. This might be enough to counter objections about generalization, since “[w]hat appear to be criticisms of generalizations in general(!), may in fact boil down to criticisms of *insufficiently precise* means of generalization.” [4, p. 5]. So in this sense, ML might cancel, reduce, or at least not increase the problem of immoral generalizations.

Similarly, the problem of data privacy pre-exists the entry of ML in the medical field [42]. Privilege bias (models being developed for diseases that disproportionately affect a certain group[38]) is a problem that exists in classical statistical studies as well [25].

Cognitive biases of all sorts (availability heuristic, anchoring effects, framing effects, tendency to see false correlations, wrong probabilistic reasoning especially with small probabilities) have been shown in decisions taken by humans [53]. In particular, availability bias and anchoring effects have been observed in medical diagnosis, with increased effects for more expert doctors who rely more on non-analytical reasoning [33]. So in general, while recognizing possible machine biases, we must remember the numerous biases shown in human reasoning for comparison.

To summarize, the trade-offs identified in medical ML (group fairness and individual fairness, privacy and predictability, fairness and predictability) exist independently of the system used to make decisions. They are not inherent to the technology used to solve them, ML, but to the goals and requirements of the system. Similarly, biases of ML systems can be tracked and compared to biases in human reasoning. ML can be used to proactively advance health equity (beneficence), and not only to avoid harms (non-maleficence) [34, 38]. This however requires deciding on a fairness measure to enforce (and how much to weigh it against performance), that will in most cases be in contrast with other fairness definitions [53].

### 4.3 Potential benefits of Machine Learning, and the sin of perfection

How can ML actively help advance health equity and fairness? Firstly, it imposes the need for precise definitions of what is meant by terms like ‘discrimination’ and ‘equity’ [53].

Secondly, it forces the developers of the system to choose precise weights for the principles that they want to respect, and explicitly accept the existence of trade-offs that are inherent to the problem. Thirdly, “it can significantly reduce one of two potential sources of bias and discrimination [...] *intrinsic* bias” [53, p. 672], by removing the influence of (unknown) prejudice and emotions from decisions (although they might still be present in the training data). Fourthly, it makes the goals and evaluation metrics (and their implied definition of what a ‘good’ solution looks like) clear. Knowing that those goals influence the results [42] (for example, pure efficiency potentially leads to the propagation of health inequities [38]), the importance of each objective has to be decided upon (and with it, the position on the Pareto front). To summarize these advantages, we can say that ML, despite the typical complaints about its inscrutability, in a way helps enforce the transparency of the decisions taken, by requiring the ethical position to be written down explicitly [52]. This is an important factor especially when comparing their use to unaided practice, where the practitioner’s values are necessarily at least indirectly influencing their decisions, probably without being stated precisely [53]. Let us reiterate that, in many cases, ML technologies would simply assist human decision-makers in existing tasks [36]. This is partly due to laws that make fully automated decisions in certain cases impossible, such as the European Union’s GDPR [53]. It must be noted that, while potentially being a solution to legal headaches, support instead of replacement does not remove the responsibility from the development chain of CDSS, since such systems might lead to the overreliance of the users on the tools and to the deskilling of practitioners [36].

In our opinion, the analysis of the impact of ML technologies must be relative to standard human decision-making (probably taking a utilitarian point of view, as is typical for stochastic problems; see [19]), and not an absolute decision about whether they act ‘perfectly fairly’ or not (which is mathematically limited in every decision system anyway [53]). A partial reason for public distrust of algorithmic solutions might be a wrong image of doctors as invulnerable and perfect figures, partly protected by physicians trying to keep a “symbolic facade of professional competence” (while privately recognizing the risks and errors of their practice) [50]. We argue that attention must be paid to always compare what is expected from algorithmic solutions and what is expected from the humans executing those tasks presently. In the case of explainability, an often-requested characteristic for ML systems in health, double standards can be shown between what is expected from ML and what practitioners currently do [53]. One can reply to the critique of ML systems as black boxes by pointing out that “the human brain, too, is largely a black box” [53, p. 666], and that humans often provide ex post explanations for decisions that are influenced by emotions and reflect mistaken rationalizations [53]. A comparison to current unaided (by ML) practice is however rarely done, and is sometimes even impossible to do because of the impossibility and/or illegality of collecting statistics about human deciders [52]. Furthermore, it might be particularly hard to do fair comparisons in the public discussions, since generally “humans are judged by their intentions, while machines are judged by their outcomes” [21, p. 139]. So, as long as fairness problems shown in human clinical decision-making are seen as unintentional, it might be difficult to argue for an effective advantage of

using algorithmic decision-making tools. Nonetheless, we hold that such comparisons are needed for a meaningful discussion on the possible use of ML systems.

The visceral resistance to the use of any technological system that shows any behavior deemed as unjust might be stopping improvements in overall care, which can be considered problematic. Do we want, for example, to refute to apply any system that does not lead to equalized outcomes but only equalized benefit [38]? How do we justify keeping the unfair status quo by avoiding solutions that would improve care in general and stratified across sensitive groups, just because those solutions do not perfectly solve the problem? Despite some resistance of doctors to evaluation (at least partially driven by fears of blame) [50], quantification of errors and fairness are needed to compare different solutions. Furthermore, the time efficiency and cost of the ML applications with respect to unaided clinicians must be taken into account. As discussed before, trade-offs are made about budget allocation and doctors' time prioritization, so reducing costs and time might allow improvements in other sectors of health care by displacing the saved resources (this way of phrasing the argument might be more positively received than simply mentioning saved costs [45]). Especially considering the benefit potential of AI in medicine, exceptionalism for the application in this field seems unjustified [17].

Since our knowledge about the world is partial and hence stochastic, we have to accept that any decision is subject to uncertainty and hence probably imperfect. In particular, the inevitability of errors in the medical field (be them active or latent) is widely recognized by practitioners themselves [50]. We hold that the advantage of actively, empirically trying an algorithmic solution, despite the uncertainty about its results, is the positive feedback loop that it creates. ML solutions will evolve over time based on the results they get, and can be monitored at deployment time (for example through failure auditing [6]). This approach would be particularly beneficial in medicine, which as a field suffers from a lacking error culture as compared to industries such as aviation and nuclear energy [50]. While the common suspect for this lack is the 'culture of blame' observed in medicine, other factors can be pointed out, such as the fear of external blame, the attempt to maintain a figure of competence, the normalization of errors, the revulsion to management, the individualistic culture, the skepticism about external non-expert observers, and collegiality [50]. These aspects would on a first analysis not influence ML systems, which should more rationally respond to errors.

When applying ML solutions, the fact that we change the way decisions are made will change the underlying data distribution and offer us more insights about the real sources of group differences (instigating some kind of population-level behavioral change health [36]). For example, pushing to correct for historical bias by applying equal allocation principles [38] will give us more diverse data based on which to infer the causes of past differences, and potentially reduce performance gaps [6]. A possible solution is hence to develop the system that we deem the more appropriate with the current knowledge, accept the imperfection and improve them over time as they get results. In a high-stakes context such as medicine, it is likely that stronger performance-based auditing (evaluating the outcomes) will be needed, and not only accreditation-based auditing (judged by experts) [53]. The close monitoring

and improvement over time of ML systems (pushed for in example in [38]) will potentially reduce the impact of dataset shift with the integration of new data. While accepting the need for strong standards on ML systems, we point out once again that those standards “*should* be applied consistently across the board, regardless of whether we are dealing with machines or humans” unless there is “some compelling political, economic or social justification to the contrary” [53, p. 678].

The optimism about the possibility of long-term improvements of health care through ML must be counterbalanced by the observation that the dynamics of the entire ecosystem make it very difficult to predict its evolution. Furthermore, since the problems that are dealt with are stochastic in nature, the developed solutions will necessarily be stochastic and entail some kind of trade-off, typical of technological innovations [19]. Dealing with stochastic problems requires a weighting of benefits and risks by their probabilities [19]. On a positive note, statistical reasoning is exactly what ML is good at (and humans are not) [52]. Unfortunately, very little work has been done in ML to assess the evolution of the data distribution when decisions are taken by ML systems adjusted for fairness. Economics literature in affirmative action may be helpful in analyzing the problem [9].

## 5 Further ethical trade-offs in medical Machine Learning systems

We have argued that most of the fairness trade-offs identified at the intersection of machine learning and medicine are not new, but rather that the preexisting ones are preserved, increased or possibly decreased. However, the deployment of machine learning methods in the medical context does introduce new trade-offs into medicine apart from the fairness domain. This might involve trade-offs that uniquely emerge from the technology of medical ML [42], so let us zoom out of the fairness domain to see what is happening when ML and medicine are combined.

### 5.1 Trade-offs

ML applications in medicine are often discussed as a human vs. machine situation - where the medical ML system outperforms the human they should and in the near future will be substituted. However, this creates a binary decision that is hard to make, especially with ML systems which can involve a good amount of uncertainty. It also creates an environment where humans are competing with a machine for the prerogative of interpretation, which contains some understanding of ML systems as some kind of autonomous systems. This autonomy is rather imaginary since as of today, medical ML systems still need their decision making process to be started and evaluated by humans. Thus, we are left with a distinction between ML systems as a tool or a machine as it is described in [52] and as it is argued, this distinction is often made based on familiarity - new developments are machines and will only be called tools when they grow older and people get used to them.

Recently, different studies found that combining ML and human evaluation can achieve better results than either of the two on their own [28, 39, 43, 46]. This suggests that human practitioners should use medical ML systems to support their decisions and increase their efficiency rather than be replaced by them, which is in line with understanding medical ML systems as tools. So instead of a binary decision we are left with a new situation that fits our understanding of trade-offs. How are ML and human evaluation best combined to achieve the best accuracy, how much human involvement do we want or need? This might heavily depend on the task at hand. For example, for skin cancer classification where the input is only a cropped image of the potential carcinoma or melanoma, the algorithms decision alone might be enough. However, for identifying diseases in a breast X-ray, a much broader task than skin cancer classification, algorithmic and human judgment might need to be combined for the optimal solution. One of the studies mentioned above found that especially for harder cases the assisted accuracy was very high compared to the unassisted accuracy when the ML model’s prediction was correct, but that it was also painfully low in cases where the ML model’s prediction was incorrect [28]. Moreover, there is also evidence that AI might especially improve the performance of less experienced practitioners like those who are still in training while those who are already experienced would not profit as much [39]. This brings another level into the trade-off because with a necessitation to use the ML system, experienced practitioners who already perform on a similar level as the system might even be hindered by another step in their workflow and thus, the overall performance could decrease. In fact, for CDSS without ML components it was already observed that more experienced doctors ignore its assistance more often without performing less good [44].

This trade-off is further complicated by the question of responsibility. Naturally, the more ML tool and physician interact the harder it gets to identify where potential mistakes come from and thus understand whether a mistake by the practitioner or an error in the medical ML system are to blame [23]). This can be problematic since many applications of medical ML systems involve high-stakes scenarios where errors should be avoided at all cost or if they happen should be eradicated as fast as possible.

A clear responsibility framework is also important for fostering trust in the medical profession which leads us to the next trade-off. For a good relationship between patient and health care practitioner, trust is of the utmost importance [10]. For the application of ML in medicine and health we can identify a multi-faceted trade-off between trust in the system and accuracy that already starts with what we just mentioned, but continues far beyond that. It involves a trade-off that is standard to ML but grows to a great importance especially in applications like medical ML systems, namely the trade-off between explainability and accuracy [27, 46]. Most medical ML systems today perform worse as soon as some kind of interpretability framework is built in, leading to the question of how important explainability is for the application [32]. While it is not yet clear, how an explainability framework does influence the work of practitioners studies have shown that it would increase trust in the system, from the practitioners as well as from the patients side [13, 46]. It seems that explainability might be a decent solution to gain the trust needed, but it could in turn also worsen the accuracy thus actually making the distrust in



the technology more reasonable. This can be related to another, rather philosophical level of this trade-off: If a person is skeptical about using ML on their diagnostic case, how can we trade-off a potentially better diagnosis against respecting the persons wish which could lead to a less accurate or even wrong diagnosis? Currently, there is no right answer to this question, it will depend on what society and the medical profession think is more important. How much do physicians need to understand the tools they are using? (Education of an AI-literate workforce [20]).

If it is impossible to reach a conclusion, for example because there is not enough data, ML tools should be transparent about that and indicate that they cannot make a decision rather than making a bad informed decision. [23]

Some scholars argue that while a mistake by a human practitioner only affects a small amount (often only one) people, a mistake by an algorithm that is deployed on many hospitals will happen more often and thus affect more people [36]. This could be seen as a trade-off between the scale of deployment of a technology and the severity of mistakes. However, this argument can also be seen as flawed since although one practitioner might not repeat a certain mistake, other practitioners not involved in this situation might well do. Thus, the only argumentation here could be that the errors are not as systematically spread as with ML tools, although even that might be an overstatement.

Another Trade-off can be seen between empathy in human practitioners and a more standardized way of tackling tasks in ML tools [36]. What do we understand as good health care, only the right diagnosis or psychological well-being during treatment, etc. as well? ethical question: should we predict death? ([46] Table 3) [20] talk about triage by ML

In particular, the use of ML as assisting systems rather than replacements of clinicians altogether complicates the discussion about biases further. The end effect of the integration of ML tools in medical practice is a complex function of the interaction of their results and their usage by clinicians on patients [38, p. 4].

## 5.2 We have seen that before

ML tools are said to be able to increase efficiency in hospitals and prevent unnecessary hospital visits, thus reducing pressure on care workers and doctors, which is certainly a good thing [23]. However, it will be important to take a holistic approach towards health care in the future. ML tools are too often seen as the holy grail to solve problems when in fact they are just tools that will not tackle structural problems without using them to do so. For example, in current health care systems reduced workload of care workers has the potential to lead to a reduction in the workforce because it is a way to save money. However, this would then not lead to an actual improvement for patients but only to a potential financial reward. This can be seen as a trade-off between monetary outcomes and spendings on the one and the patients experience and care on the other hand. While this is an issue that is already existing, ML tools bring another perspective to it since they have the potential to increase as well as heavily decrease the patients experience in hospitals.

Often, ML tools only work for specific tools, i.e. detecting one or a couple of diseases

in an X-ray. While the accuracy rate here is often high the broadness of the analysis is very limited compared to a doctor [46]. This could be identified as a trade-off between high accuracy with a narrow focus on the one hand and lower accuracy with a broader focus on the other.

The current way of handling medical data differs heavily from the way data is used in ML [20]. Unfortunately, to make ML tools work properly there is a need for huge amounts of data that will be shared with the respective companies and researchers. This creates a trade-off between the classical handling of medical data and a necessary data collection.

Health care systems around the world are more or less privatized, depending on the country. However, in the case of ML tools a lot of research and development is driven by big companies like Alphabet or IBM [36]. This makes sense since those companies are driving ML research in general but it poses the question whether we want to give such an important issue completely out of public and into private hands. While the privatization of health care was already posing problems before ML tools and they are in fact seen as a solution for the existing problems [36, 46] the questioning of privacy and trust is increased as well. Thus, this can be seen as a trade-off between the speed of development - arguably, big tech companies will be fast in bringing ML tools to the market - and privacy and trust issues.

The development and deployment of ML tools in medicine and health care will and does already cost a lot of money [20]. At the same time, the health care system in general in countries like the US is heavily underfunded. So much so that life expectancy began to decrease again in the US [46]. Thus, there can be a trade-off identified between the financing of ML tools and the health care system in general. If the huge investments in ML tools will only benefit a small wealthier part of society, those investments are questionable if the health care systems continue to be underfunded. This is even more the case since there are not yet many ML tools ready for clinical application which makes this money an investment into the future while there persist acute issues that would need to be tackled here and now.

WTF? The savings would come from a combination of deployments: lower medical costs and reduced losses from low productivity and sick day ([23] page 148)

Another trade-off exists between the way medical devices are traditionally approved for (clinical) applications and how software is usually deployed and constantly updated [20].. While this problem might also exist with software that is already deployed in other ways in medicine and health care, ML tools take it to a new level. Here, updates might involve newly trained algorithms with a new data background which might have achieved different performance benchmarks. How should this agile updating be weighed against traditional and more accurate, but slower ways of approving tools for clinical application? In the US, the FDA already reacted by creating easier paths for improvement for this kind of software but the success of this pathway is still indeterminate. Today, regulation processes are often such that the model is locked in place before deployment. This makes it easier to regulate them but misses out on their potential to learn and increase functionality on the fly.

How do we study the clinical efficacy of ML tools? Is a randomized controlled trial ethical? Because with normal medical trials we do not have an alternative working treat-

ment, we just compare it with nothing. Thus, we do not withhold something from patients. However, if ML tools (wrongfully) decide against treatment we actively withhold treatment from patients which might in extreme cases lead to their deaths. (Grote und Genin)

Many of the trade-offs discussed can essentially be broken down to one question: How much do we benefit from the use of ML in medicine? What might be bad for us, for example could the digitalization and sharing of our health data lead to misuse by health care providers? If I know that I will most probably benefit from sharing my data on the other hand, I will be more likely to do so. This would create a very general trade-off between benefits and caveats of ML tools in medicine and health care. ([46] Increased Efficiencies) ([20] Transparency) Can digital health care be for everybody? what about people who do not have digital devices or don't want to use them?

But this also leads to the question whether these trade-offs are actually new or - as discussed for fairness trade-offs above - if they are just new editions of trade-offs humanity has seen before in either ML, medicine and health care or technology in general.

This is certainly not only a problem in medicine, but medicine has always been a discipline where trust is particularly important. Machine Learning will bring a new twist to this issue and medicine is in that sense a unique challenge for Machine Learning.

## 6 Conclusions

We were able to show that trade-offs are omnipresent in technology, that they usually cannot be avoided and that they can be formalized as Pareto optimality. Three main trade-offs involving fairness were identified in ML: The choice of the fairness measure which involves the decision between individual and group fairness, the trade-off between accuracy and fairness which results from the fact that ensuring fairness often leads to a lower prediction accuracy, and the trade-off between privacy and accuracy, where ensuring privacy has the same effect as ensuring fairness.

A number of trade-offs can be identified in health care, at different levels, for example in policy, treatment, or research. Often, trade-off situations are described in terms of moral dilemmas. We argued that they are fundamentally the same thing, but people tend to reject trade-off talk when sacred values (such as human life) are involved [45]. Fairness issues have been identified in the current practice of medicine, even when unaided by algorithmic tools. This is evident also from the fact that so-called biased data, partly responsible for algorithmic fairness issues, is generated in those situations. While we should take inspiration from biomedical ethics literature, since fairness analysis is largely domain-dependent, we must at the same time recognize that principled approaches might lead to disagreement rather than solve it.

Many fairness-related issues identified in the fair-ML literature in medical applications are not strictly caused by ML. Rather, they are made more explicit by requiring mathematical formulations. We hold that ML tools should be analyzed with respect to fairness relatively to human deciders (doctors), rather than with binary fair/not fair questions. Errors and imperfections have to be expected, and are unavoidable in stochastic situations. In

this respect, care must be taken in recognizing double standards for humans and machines when present, keeping in mind the limitations of medicine as is practiced today. ML can be used to actively improve fairness.

## References

- [1] Christopher Alexander. *Notes on the Synthesis of Form*. Vol. 5. Harvard University Press, 1964.
- [2] Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. “Differential privacy has disparate impact on model accuracy”. In: *Advances in Neural Information Processing Systems* 32 (2019).
- [3] Richard Berk et al. “Fairness in criminal justice risk assessments: The state of the art”. In: *Sociological Methods & Research* 50.1 (2021), pp. 3–44.
- [4] Reuben Binns. “Fairness in machine learning: Lessons from political philosophy”. In: *Conference on Fairness, Accountability and Transparency*. PMLR. 2018, pp. 149–159.
- [5] Reuben Binns. “On the apparent conflict between individual and group fairness”. In: *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 2020, pp. 514–524.
- [6] Richard J Chen et al. “Algorithm fairness in AI for medicine and healthcare”. In: *arXiv preprint arXiv:2110.00603* (2021).
- [7] Andrew Chester et al. “Balancing utility and fairness against privacy in medical data”. In: *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE. 2020, pp. 1226–1233.
- [8] Alexandra Chouldechova. “Fair prediction with disparate impact: A study of bias in recidivism prediction instruments”. In: *Big data* 5.2 (2017), pp. 153–163.
- [9] Alexandra Chouldechova and Aaron Roth. “A snapshot of the frontiers of fairness in machine learning”. In: *Communications of the ACM* 63.5 (2020), pp. 82–89.
- [10] Chalmers C Clark. “Trust in medicine”. In: *The Journal of medicine and philosophy* 27.1 (2002), pp. 11–29.
- [11] Rachel Cummings et al. “On the compatibility of privacy and fairness”. In: *Adjunct Publication of the 27th Conference on User Modeling, Adaptation and Personalization*. 2019, pp. 309–315.
- [12] Francois Dionne and Craig Mitton. *Health Care Trade-Offs: A Necessary Reality For Every Health System*. Mar. 20, 2018. URL: <https://www.healthaffairs.org/doi/10.1377/forefront.20180316.120106%7D> (visited on 02/26/2022).
- [13] William K Diprose et al. “Physician understanding, explainability, and trust in a hypothetical machine learning risk calculator”. In: *Journal of the American Medical Informatics Association* 27.4 (2020), pp. 592–600.
- [14] John C Duchi, Michael I Jordan, and Martin J Wainwright. “Privacy aware learning”. In: *Journal of the ACM (JACM)* 61.6 (2014), pp. 1–57.
- [15] Cynthia Dwork et al. “Fairness through awareness”. In: *Proceedings of the 3rd innovations in theoretical computer science conference*. 2012, pp. 214–226.

- [16] Robert Eiss. “Confusion over Europe’s data-protection law is stalling scientific progress”. In: *Nature* 584.7822 (2020), pp. 498–499.
- [17] Richard Ribón Fletcher, Audace Nakeshimana, and Olusubomi Olubeko. “Addressing fairness, bias, and appropriate use of artificial intelligence and machine learning in global health”. In: *Frontiers in Artificial Intelligence* 3 (2021), p. 116.
- [18] Pratyush Garg, John Villasenor, and Virginia Foggo. “Fairness metrics: A comparative analysis”. In: *2020 IEEE International Conference on Big Data (Big Data)*. IEEE. 2020, pp. 3662–3666.
- [19] Russell Hardin. “Ethics and stochastic processes”. In: *Social Philosophy and Policy* 7.1 (1989), pp. 69–80.
- [20] Jianxing He et al. “The practical implementation of artificial intelligence technologies in medicine”. In: *Nature medicine* 25.1 (2019), pp. 30–36.
- [21] César A Hidalgo et al. “Moral functions”. In: *How humans judge machines*. MIT Press, 2021, pp. 123–147.
- [22] Vincent KY Ho. “Medicine, methodology, and values: trade-offs in clinical science and practice”. In: *Perspectives in Biology and Medicine* 54.2 (2011), pp. 243–255.
- [23] Denis Horgan et al. “Artificial intelligence: power for civilisation—and for better health-care”. In: *Public health genomics* 22.5-6 (2019), pp. 145–161.
- [24] C-L Hwang and Abu Syed Md Masud. *Multiple objective decision making—methods and applications: a state-of-the-art survey*. Vol. 164. Springer Science & Business Media, 2012.
- [25] Gabrielle Jackson. *The female problem: how male bias in medical trials ruined women’s health*. Nov. 13, 2019. URL: <https://www.theguardian.com/lifeandstyle/2019/nov/13/the-female-problem-male-bias-in-medical-trials> (visited on 02/02/2022).
- [26] Faisal Kamiran, Toon Calders, and Mykola Pechenizkiy. “Discrimination aware decision tree learning”. In: *2010 IEEE International Conference on Data Mining*. IEEE. 2010, pp. 869–874.
- [27] Christopher J Kelly et al. “Key challenges for delivering clinical impact with artificial intelligence”. In: *BMC medicine* 17.1 (2019), pp. 1–9.
- [28] Amirhossein Kiani et al. “Impact of a deep learning assistant on the histopathologic classification of liver cancer”. In: *NPJ digital medicine* 3.1 (2020), pp. 1–8.
- [29] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. “Inherent trade-offs in the fair determination of risk scores”. In: *arXiv preprint arXiv:1609.05807* (2016).
- [30] John Launer. “Medicine and the art of trade-offs”. In: *Postgraduate Medical Journal* 96.1139 (2020), pp. 575–576. ISSN: 0032-5473. DOI: 10.1136/postgradmedj-2020-138575. eprint: <https://pmj.bmj.com/content/96/1139/575.full.pdf>. URL: <https://pmj.bmj.com/content/96/1139/575>.

- [31] Suyun Liu and Luis Nunes Vicente. “Accuracy and fairness trade-offs in machine learning: A stochastic multi-objective approach”. In: *arXiv preprint arXiv:2008.01132* (2020).
- [32] Yi Luo et al. “Balancing accuracy and interpretability of machine learning approaches for radiation treatment outcomes modeling”. In: *BJR— Open* 1.1 (2019), p. 20190021.
- [33] SÝlvia Mamede et al. “Effect of availability bias and reflective reasoning on diagnostic accuracy among internal medicine residents”. In: *Jama* 304.11 (2010), pp. 1198–1203.
- [34] Melissa D McCradden et al. “Ethical limitations of algorithmic fairness solutions in health care machine learning”. In: *The Lancet Digital Health* 2.5 (2020), e221–e223.
- [35] Aditya Krishna Menon and Robert C Williamson. “The cost of fairness in binary classification”. In: *Conference on Fairness, Accountability and Transparency*. PMLR. 2018, pp. 107–118.
- [36] Jessica Morley et al. “The ethics of AI in health care: A mapping review”. In: *Social Science & Medicine* 260 (2020). ISSN: 0277-9536. DOI: <https://doi.org/10.1016/j.socscimed.2020.113172>. URL: <https://www.sciencedirect.com/science/article/pii/S0277953620303919>.
- [37] Bryan Oronsky et al. “Medical Machiavellianism: the tradeoff between benefit and harm with targeted chemotherapy”. In: *Oncotarget* 7.8 (2016), p. 9041.
- [38] Alvin Rajkomar et al. “Ensuring fairness in machine learning to advance health equity”. In: *Annals of internal medicine* 169.12 (2018), pp. 866–872.
- [39] Pranav Rajpurkar et al. “AI in health and medicine”. In: *Nature Medicine* (2022), pp. 1–8.
- [40] Will Rinehart and Allison Edwards. “Understanding job loss predictions from artificial intelligence”. In: *American Action Forum. Org.* 2019.
- [41] M Teresa Ruiz and Lois M Verbrugge. “A two way view of gender bias in medicine.” In: *Journal of epidemiology and community health* 51.2 (1997), pp. 106–109.
- [42] Marianne WMC Six Dijkstra et al. “Ethical considerations of using machine learning for decision support in occupational health: An example involving periodic workers’ health assessments”. In: *Journal of Occupational Rehabilitation* 30 (2020), pp. 343–353.
- [43] David F Steiner et al. “Impact of deep learning assistance on the histopathologic review of lymph nodes for metastatic breast cancer”. In: *The American journal of surgical pathology* 42.12 (2018), p. 1636.
- [44] Reed T Sutton et al. “An overview of clinical decision support systems: benefits, risks, and strategies for success”. In: *NPJ digital medicine* 3.1 (2020), pp. 1–10.
- [45] Philip E Tetlock. “Thinking the unthinkable: Sacred values and taboo cognitions”. In: *Trends in cognitive sciences* 7.7 (2003), pp. 320–324.

- [46] Eric J Topol. “High-performance medicine: the convergence of human and artificial intelligence”. In: *Nature medicine* 25.1 (2019), pp. 44–56.
- [47] Stephen Toulmin. “How medicine saved the life of ethics”. In: *Perspectives in biology and medicine* 25.4 (1982), pp. 736–750.
- [48] Jeroen Van Den Hoven. “Information technology, privacy, and the protection of personal data”. In: *Information technology and moral philosophy* 301 (2008).
- [49] Irene Van Staveren. *The ethics of efficiency*. Tech. rep. 2007.
- [50] Justin J Waring. “Beyond blame: cultural barriers to medical incident reporting”. In: *Social science & medicine* 60.9 (2005), pp. 1927–1935.
- [51] Susan Wei and Marc Niethammer. “The fairness-accuracy Pareto front”. In: *Statistical Analysis and Data Mining: The ASA Data Science Journal* (2020).
- [52] Robert C Williamson. “The AI of Ethics”. In: *Machines We Trust: Perspectives on Dependable AI*. Ed. by Marcello Pelillo and Teresa Scantamburlo. MIT Press, 2021. Chap. 9, pp. 139–160.
- [53] John Zerilli et al. “Transparency in algorithmic and human decision-making: is there a double standard?” In: *Philosophy & Technology* 32.4 (2019), pp. 661–683.
- [54] Indre Zliobaite. “On the relation between accuracy and fairness in binary classification”. In: *arXiv preprint arXiv:1505.05723* (2015).