# Fair Decision-making Under Uncertainty

Wenbin Zhang and Jeremy C. Weiss
Carnegie Mellon University, United States
Email:{wenbinzhang, jeremyweiss}cmu.edu

*Abstract*—There has been concern within the artificial intelligence (AI) community and the broader society regarding the potential lack of fairness of AI-based decision-making systems. Surprisingly, there is little work quantifying and guaranteeing fairness in the presence of uncertainty which is prevalent in many socially sensitive applications, ranging from marketing analytics to actuarial analysis and recidivism prediction instruments. To this end, we study a longitudinal censored learning problem subject to fairness constraints, where we require that algorithmic decisions made do not affect certain individuals or social groups negatively in the presence of uncertainty on class label due to censorship. We argue that this formulation has a broader applicability to practical scenarios concerning fairness. We show how the newly devised fairness notions involving censored information and the general framework for fair predictions in the presence of censorship allow us to measure and mitigate discrimination under uncertainty that bridges the gap with real-world applications. Empirical evaluations on real-world discriminated datasets with censorship demonstrate the practicality of our approach.

*Index Terms*—Fairness, censorship, survival analysis

## I. INTRODUCTION

The role of artificial intelligence is expanding and is transforming many walks of life, such as the screening of job applications, the setting of insurance rates, the targeting of advertising, the allocation of health resource and the approval of mortgage loans [1]–[3]. This trend is likely to continue, and the performance of the AI systems may match or even surpass human level [4]. However, as AI permeates our lives in domains of high societal impact, there is increased concern regarding the fairness and accountability of AI-based decision-making systems, which have been voiced within and beyond the AI community [5]–[7]. In a recent and widely popular investigation conducted at ProPublica, authors Chen and Wong analyzed state-of-the-art clinical prediction models, concluding that these models are biased against black patients by systematically underperforming on them even when the treatment is aimed at a particular type of cancer that disproportionately impacts them [8]. Because AI-based decision-making can be as biased as human and can even exacerbate disparity, there is therefore an urgent need to consider fairness in AI algorithms in order to maximize AI benefits for social good.

This has led to an active area of research into quantifying and mitigating AI unfairness for the sake of providing fairness-aware decision-making systems, *i.e.*, systems that are not unduly biased for or against certain individuals or social groups. Most existing fairness notions seek to evaluate the bias of a decision-making process by means of two particular forms, *disparate impact* and *disparate treatment*, in legal texts [9].

While the former [10] asks for approximate parity when allocating a more favorable outcome across different demographic groups defined by sensitive attributes (*e.g.*, race or gender), the latter [11] guarantees similar individuals are treated similarly irrespective of sensitive characteristics. Moreover, most of the work in the existing literature tackles the fairness problem by assuming the presence of class label, in which the fairness notions are defined based on the class label either actual or predicted, and the same predictive model is trained contingent upon for new instance prediction [12], [13]. There is limited work focuses on *censoring settings* in which the class label could be absent due to censoring on the time to an event of interest, preventing existing fairness notions and approaches from being applied.

However, the censoring phenomenon widely exists in many real-world applications. In addition to the aforementioned discriminatory clinical prediction study in which the patient's true time to relapse or hospital discharge could be unknown for various reasons, censoring is common when profiling customers for business planning in marketing analytics, predicting criminal recidivism for recidivism prediction instruments, to name a few. Thus, it is necessary to design decision-making systems free from discrimination while considering the characteristics of censoring settings, which is underexplored and brings unique challenges: **i) The lack of fairness definitions under uncertainty.** Although more than twenty measures of fairness have been proposed in the literature, formalizing fairness is a hard topic per se and most of them necessitate the certainty on the class label which is infeasible in the censoring settings [14]. There is no existing fairness definition explicitly involving censorship information to quantify fairness under uncertainty. **ii) Addressing fairness in the presence of censoring.** Different from the supervised fairness methods assuming the availability of class label either actual or predicted by-design, extra care must be taken to ensure the automated decision-making system is independent of sensitive attributes-based harmful stereotypes in the presence of censoring. Addressing both unfairness and censored distribution simultaneously is challenging as the censorship accompanies and complicates the biased decision regions. In addition, directly combining/transferring from respective domains is not straightforward, sophisticated design is therefore warranted. **iii) Theoretically navigating the interplay between fairness and censoring.** Existing fairness-aware studies have largely been focused on addressing and analyzing fairness in the presence of class label. However, when there is an uncertainty on the class label, existing fairness theoretical analyses cannot

apply. In addition, such a theoretical analysis is also desired so as to allow for practitioners and policy makers to navigate and customize according to their business objectives [5].

To the best of our knowledge, this is the first work that can address the above challenges and provides a generic framework that specifically incorporates censoring information for fair decision-making. More specially, the novelty of this research comes from four aspects:

- We formulate a new problem of fair decision-making under uncertainty. Then, we devise corresponding fairness notions to measure unequal treatments in the presence of censorship and attribute with different semantic meanings, thus providing necessary complements to existing fairness notions in the literature.
- A respective fairness-aware learner is developed for learning with censored data that are common in many real-world applications. The proposed learner specifically accounts for censored information in the model building so as to ensure accurate predictions while minimizing unfairness in censoring settings.
- A theoretical analysis on the interplay between fairness and censorship is conducted, which broadens the existing fairness theory to new types. Such an analysis enables model-agnostic evaluation of fairness and censorship interplay and is also of practical purpose to help practitioners navigate through their business needs.
- Quantitative and qualitative experiments on a set of real-world discriminated datasets demonstrate that the proposed fairness notations and debiasing algorithm are indeed capable of measuring and guaranteeing fairness in the presence of uncertainty.

The remainder of this paper is organized as follows. We will start with a review of related work on fairness-aware learning and survival analysis in Section II. Next, we will discuss the problem definition along with notations used in this work in Section III. In Section IV we propose, to the best of our knowledge, the first fairness notions and debiasing algorithm that explicitly involve censoring information for fairness under uncertainty. Our experimental methodology and results are summarized in Section V. We end with our conclusions in Section VI.

## II. RELATED WORK

Relevant to our work is the work that quantifies and mitigates machine learning algorithmic unfairness and work that handles with censored data, in particular analyzing censored data with decision tree models.

### A. Quantifying Unfairness

While machine learning increasingly permeate facets of life, significant concerns on the unfair and discriminatory manner of ML-based systems have been voiced and observed [15]. The machine learning community has responded by proposing a growing body of fairness notions to measure the level of discrimination along with a number of approaches to mitigate

bias in order to provide fairness-aware decision-making systems [16].

The broad set of existing mathematical formulations of fairness can be typically divided into two main families, individual fairness and group fairness, based on whether they evaluate fairness at the individual level or the group level, respectively [17]. The former definitions aim to ensure that similarly situated individuals are treated similarly. The seminal work is the notion proposed by Dwork et al. [11], which requires similarly situated individuals to receive similar probability distributions over class labels to prevent unequal treatments. One of the prerequisites of this line of work is the demanding of a task-specific similarity metric that is suitable to measure how similar two individuals are depending on the task at hand. In practice, however, such a proper similarity metric cannot be trivially specified or it is even not possible to do so. Take the individuals with censorship as an example, a proper similarity metric is hard to be specified unless disregard the censored information which is of great importance and cannot simply be ignored. The lack of task-specific similarity metric has therefore been a major obstacle for the wider adoption of individual fairness [18].

On the other hand, group fairness notations normally first specify a set of subgroups defined by sensitive attributes (*e.g.*, race or gender), then preserve fairness by asking for group level approximate parity of some statistic over class labels either predicted or actual. For example, statistical or demographic parity, the representative formulation belonging to this category, measures whether the desired outcomes are equally distributed across different demographic groups [19]. Compared to the dependence of predicted class label of statistical parity, equalized odds [20], another representative notation, further depends on the actual class label by considering the true classification difference between different subgroup.

Although many fairness notions exist, most of them, as previously discussed, formulate fairness depend on class label either actual or predicted, thus limiting their applicability in censoring settings which is prevalent in real-world applications. Keya et al. [21] directly extend the existing fairness notions to the application with censoring problems, which is the single relevant effort to the best of our knowledge. However, their definitions exclude the censoring information when measuring discrimination, which could introduce substantial bias, as censored information might be of importance and cannot simply be ignored [22]. In addition, a similarity metric, as one of the vital components, is required to be specified with the Euclidean distance employed in their study, which could lead to inappropriate similarity evaluation due to censorship as well as the potential distinct semantic meanings of different attributes [23]. Our proposed fairness notions alleviate such limitations by explicitly involving censoring information and is free of similarity metric specification.

### B. Mitigating Unfairness

The fairness definitions above could be directly used or slightly modified as a constraint or a regularizer to enforce

fairness, leading to three categories of mechanisms to guarantee fairness: i) pre-processing approaches, ii) in-processing approaches and ii) post-processing approaches, depending on the intervention occurs at the input/data layer, the algorithm design or the output/results of the model, respectively.

The first strategy, *pre-processing approaches*, deals with bias in the input/data level. The underlying assumption is that in order to learn a fair classifier, the training data should be discrimination-free. Therefore, the methods in this category, try to "correct" the data to ensure fairness in the representation of different communities. This comprises a popular category of methods as it is model-agnostic and therefore, can be employed by any classifier. For instance, massaging [24] changes the data distribution by re-labeling some of the instances in order to neutralize discriminatory effects. The instance to be altered are those close to the decision boundary, which is provided by a ranker. The method is applicable to binary classification problems, so the re-labeling is actually label swap. Massaging might result in reverse discrimination, *i.e.*, communities that were favored are now discriminated. In [25] an extension of massaging is proposed that also deals with the problem of reverse discrimination. Reweighting [26], on the other hand, assigns different weights to the different communities, *e.g.*, the deprived community will receive a higher score comparing to the favored community. However, the underlying assumption of such type of approaches does not necessarily hold as inherit bias could still persist due to proxies/correlations between features.

The second category, *in-processing techniques*, directly modifies the learning algorithm to ensure that it will produce non-discriminating results. Typically, these methods are algorithm-specific. As an example, the fairness gain, reflected by the reduction in discrimination, is incorporated into the splitting criterion for fair tree induction [27]. This model was later extended as an ensemble approach offering additional adaptability and flexibility of fairness [28]. In [29], the Mann Whitney U statistic, a rank-based non-parametric independence test, is leveraged to measure the correlations between class label and sensitive attribute, then further reformulates it as a new non-convex optimization problem to mitigate the inherent bias of the data. The common assumption hold by these methods is the certainty on the class label. On the contrary, research on fairness under uncertainty has still been limited due to the significant technical challenges when addressing unfairness and censoring distribution simultaneously [21]. Our work seeks to fill this gap by jointly addressing bias reduction and censoring management.

The last category, *post-processing solutions*, "corrects" the results of a model by modifying its decision regions to ensure a fairer representation of different communities. For example, in [30], the prediction thresholds are adjusted to decrease discrimination while minimizing the effect on predictive accuracy whereas in [31], the authors shift the decision boundary of the deprived group to work against discrimination. In this category, one could also place method for building human-interpretable models [32]. We emphasize that transferring such techniques to censoring settings is not straightforward as the decision boundaries could be censored themselves due to the censored distributions in censoring settings.

## C. Survival Analysis

The critical challenge of the main outcome under assessment could be unknown for a portion of the study group, deemed censoring, hinders the use of many methods of analysis. This motivates the study of survival analysis to address the problems of partially survival information access from the study cohort [14], [33], [34].

With the ubiquitous of censored data, survival analysis has gained its popularity in applications beyond its originated medical domain ranging from customer analytics and actuaries to predictive maintenance in mechanical operations [35]–[37]. Among the various methods proposed for modeling censored data, the Cox proportional hazards model (CPH) [38] is the most commonly used in which the multiplicative relation between the risk, as expressed by the hazard function and covariates is described. Under the assumption of proportional hazards, deep neural networks have recently also been employed to better encode the nonlinearity of censored data [39], [40]. In contrast, another line of effort is the tree based methods [41], [42], particularity random forests due to its superior capabilities in handling nonlinear effect of variables and is free of restrictive assumptions such as proportional hazards [36], [43]. A comprehensive literature survey covering recent censored data modeling effort is provided in [44].

Amid the popularity of survival models, same as other ML approaches, care must be taken to ensure the fairness of these models. Our work situates in this under-explored research direction to tackle fairness in the presence of censoring. Our in-processing approach incorporates a fair splitting criterion that explicitly considers censoring information into the algorithm design to guide an accuracy-driven and fairness-oriented learning procedure. Relevantly, the Cox model is modified to ensure fair risk predictions in [21]. Three key differences are that our model: i) does not necessitate a distance metric to be specified for the wide adoption of our method, ii) explicitly includes survival information and survival time to mitigate bias in the censoring settings, and iii) is free of hyperparameter tuning to decrease the computational requirements in practice.

## III. NOTATIONS AND PROBLEM DEFINITION

In the typical fairness-aware learning settings, the discriminated data $X$ normally consists of a sequence of feature represented instances $x_1, x_2, \cdots, x_n$. Among the feature representation, a special attribute $G$ is referred as the *sensitive attribute* and its attribute values distinguish the discriminated community, *i.e.*, the deprived group, from the privileged community, *i.e.*, the favored group. In addition, instances are also described by their corresponding class labels $y_1, y_2, \cdots, y_n$. However, class labels are inaccessible in the presence of censorship.

The discriminated and censored data, in contrast to the typical data representation, therefore further contains the survival

time $T$ and an event indicator $\delta$ in addition to the observed features $x$, typically represented in the form of $(x, T, \delta)$. If the event of interest has occurred, $T$ is the actual time from the individual entered the study till the time of the event occurring, and $\delta$ becomes 1 indicating certainty on the event observation; otherwise $T$ corresponds to the elapsed time between individual entered the study and last follow-up with the individual, and the event indicator $\delta = 0$, *i.e.*, the survival time is censored [45].

When solely focus on the survival data without fairness constraints, the data is generally considered and modeled in terms of two quantitative terms, namely the hazard function and the survival function. The former models the instantaneous rate of event occurs at a specified time $t$ condition on surviving to $t$:

$$h(t|x) = \lim_{\triangle t \to 0} \frac{Pr(t < T < t + \triangle t | T \geq t, x)}{\triangle t} \quad (1)$$

While the latter is the probability that the event does not occur up to time $t$, and can be determined from the hazard function and vice versa:

$$S(t|x) = exp(-H(t|x)) \quad (2)$$

$$H(t|x) = \int_0^t h(t|x)dt \quad (3)$$

where $H(t|x)$, known as the cumulative hazard function, is the intermediate function between the hazard function and the survival function, and can be naturally interpreted as the expected number of events of interest [46].

Compared with fairness-aware learning in supervised settings, addressing discrimination bias in censoring settings leads to uncertainty on $y_1, y_2, \cdots, y_n$ which limits the applicability of the existing fairness notions. In addition, the uncertainty on $y_1, y_2, \cdots, y_n$ could also further accompany and complicate the biased decision regions. Given the discriminated and censored data $X$, the aim of fairness-aware learning under uncertainty is then to model a fair survival function $H(\cdot)$ which makes accurate predictions based on $X$ but also does not discriminate with respect to $G$ for the discriminated and censored datasets.

## IV. METHODOLOGY

This section first introduces the first of its kind fairness definitions specifically accounts for censoring, then the fair splitting criterion for the discrimination-aware and censorship information involved tree induction is introduced followed by a corresponding random forests based learning algorithm for fair risk prediction. Last, the interplay between fairness and censoring is theoretically analyzed.

### A. Censored Fairness Metrics

The presence of censorship in data limits the applicability of commonly used fairness definitions introduced in the existing fairness-aware studies. To fill this gap, we introduce two fairness metrics specifically account for model unfairness in the present of censorship to help with practitioners and policy makers arriving to fair decisions.

**Concordance imparity**. Motivated by the previously discussed ProPublica investigation on clinical prediction models [8], we first propose the *concordance imparity (CI)* to measure whether the model under consideration has systematically underperformed on certain population. Formally, we define the CI, mathematically represent as $C_\triangle$, as the largest deviation of discriminative abilities across different demographic groups of the model:

$$C_\triangle = \max_{g, g' \in G \ \& \ g' \neq g} |F(g) - F(g')| \quad (4)$$

where F refers to the *concordance fraction* evaluating the group-wise correct pairwise ordering based on its respective group members and is formulated as below:

$$F(G = g) = \frac{1}{|M(G = g)|} \sum_{i|G(x_i)=g} \sum_{j \neq i} \mathbb{1}[r(x_i) > r(x_j)] \quad (5)$$

where $r(x_i)$ is the risk score assigned to the individual $x_i$ by the model, and $M$ stands for the *permissible pair* whose shorter survival time is observed (the *impermissible pair* is denoted as $N$ to be further discussed in Section IV-D). The $M$ is mathematically defined as:

$$M(G = g) = \{G(x_i) = g, x_j \neq x_i | \delta_{min(t_i, t_j)} = 1\} \quad (6)$$

Specifically, CI first considers individual level pairwise comparison between model prediction and actual outcomes, *i.e.*, Equation (5) and (6), then measures, at the group level, *i.e.*, Equation (4) and (5), whether the discriminative ability of the model is fairly distributed across groups. The lower the concordance imparity score the fairer the model. Note that different from the previous definitions [21], the survival time and censorship information are explicitly involved in CI to avoid important information loss and introducing substantial bias.

**Fair calibration**. The concordance imparity depends on the risk score to measure the distinct discriminative capabilities on different demographic groups of the model. Here, we further introduce *fair calibration (FC)* to evaluate probability values based prediction disparate. In contrast to the risk score based prediction disparate, which is meaningful within the context of other patients' risk scores, probability values are labels for individual patients with semantic content, thus providing stakeholders and practitioners with an additional navigation for fair decisions.

In practice, FC starts with sorting the predicted probabilities at a specific time $t$ condition on surviving to $t$, *i.e.*, $S(t|x)$, for each individual demographic group as defined by $G$. Within each demographic group, the sorted predicted probabilities are further group into deciles, *i.e.*, $D = 10$ number of bins. Suppose there are 100 individuals for one certain demographic group, *i.e.*, $G = g$; each bin would contain 10 individuals

belong to this demographic group with their predicted probabilities in ascending order, so are other demographic groups. The final outcome of FC, denoted $C_\parallel$, is then evaluated as:

$$C_\parallel = \begin{cases} \text{fair calibrated,} & p(L_g(S(t|x))) \geq 0.05, \forall \text{ } g \in G \\ \text{not fair calibrated,} & \text{otherwise} \end{cases} \quad (7)$$

where $p(L_g(S(t|x))$ is the p-value of the Hosmer-Lemeshow goodness-of-fit test [47] for demographic group $g$ at a specific time $t$:

$$L_g(S(t|x)) = \sum_{i=1}^{D} \frac{(O_{ig} - n_{ig}\bar{p}_{ig})^2}{n_{ig}\bar{p}_{ig}(1 - \bar{p}_{ig})} \quad (8)$$

where $O_{ig}$, $n_{ig}$ and $\bar{p}_{ig}$ are the number of observed events, the number of individuals and the average predicted probability of bin $i$ pertaining to demographic group $g$ at time $t$, respectively. However, the exact number of observed events, $O_{ig}$, could be unobservable when bin $i$ contains individuals censored before time $t$. To this end, the D'Agostino-Nam translation [48] is employed to incorporate censoring which uses the Kaplan-Meier [49] curve estimate of bin $i$ pertaining to $g$, denoted $K_{ig}$, in place of $O_{ig}$:

$$O_{ig} = n_{ig}(1 - K_{ig}(t)) \quad (9)$$

The $L_g(S(t|x))$ follows a $\chi^2_{B-1}$ distribution, and the model is fair calibrated on the condition that all p-value pertaining to each subgroup pass the test, *i.e.*, a p-value greater than 0.05; otherwise, the model's predicted probability is deemed to be biased against those being underperformed, *i.e.*, the model's predicted probabilities only agree with favored communities' observed event rates or frequencies of the outcome. Note that the same as concordance imparity, fair calibration explicitly involves survival time and censoring information so as to quantify unfairness in the presence of uncertainty effectively.

### B. Universal Survival Difference

Random forests construct an array of base learners to increase predictive ability. It has also been extended to handle censored data, and enjoys the merits of automatically addressing the difficulties of restrictive parametric assumptions of other approaches and is capable of modeling nonlinear interactions [44]. However, existing survival random forests aim to optimize for predictive performance and does not take fairness, which we desire to add, into consideration [43]. To alleviate this limitation, we propose *Survival Universal Random Forests (SURF)* to jointly consider fair data encoding and discrimination reduction for fair forests induction by: i) introducing a new accuracy-driven and fairness-oriented splitting criterion to select the potential fair splitting candidates in the presence of censorship, ii) illustrating the way to enable the construction of the new fair random forests under censorship (c.f. Section IV-C).

In the absence of censorship, the information gain and Gini impurity are normally used to measure the certainty on

class labels during the induction of the tree for classification performance [50]. When censorship is present, which brings uncertainty on the class labels, the survival difference between splitting nodes can be instead used to measure the impurity. To explicitly involve survival time and censoring information, we use the *logrank test* [51] to distinguish the survival difference, denoted $S_\triangle$, between different groups:

$$S_\triangle = \frac{\sum_{j=1}^{k}(O_j - E_j)}{\sqrt{\sum_{j=1}^{k} V_j}} \sim N(0, 1), \quad (10)$$

where $O_j$, $E_j$ and $V_j$ are the observed number of events, the expected number of events and variance of the observed number of events, respectively. Splitting on candidates with a larger logrank test leads to similar survival intra child nodes and dissimilar survival inter child nodes.

We then propose a new fair splitting criterion that takes both predictive performance and fairness into consideration. We define the conjunctive criterion *Universal Survival Difference (USD)*, denoted $S_\vee$, as:

$$S_\vee = |\log S_\triangle(X_v, X_{v'})| - \sum_{v \in dom(f)} |\log S_\triangle(X_v^g, X_v^{g'})| \quad (11)$$

where $X_v$ (or $X_{v'}$), $v$ (or $v') \in dom(f)$ are the partitions induced by feature $f$ with $X_v^g$ and $X_v^{g'}$ represent the sub-partitions further distinguished by the sensitive attribute $G$ within the partition $X_v$. In practice, the minuend focuses on the overall survival difference between different resulting tree nodes if induced by feature $f$, while the subtrahend pays attention to the survival difference between different demographic groups within each resulting tree node.

Intuitively, USD jointly considers the predictive performance by maximizing the internode-wise overall survival difference, *i.e.*, a large value of the minuend, as well as unfairness reduction by minimizing the intranode-wise survival difference among different demographic subgroups, *i.e.*, the smaller subtrahend value, thus motivating fair split by giving priority to splitting candidates that ensures predictive performance while preventing from systematically underperforming on certain population simultaneously. The use of $log$ form in USD is for smoothing. In addition, such a combination of multiple factors is also advantageous when factors are not in the same scale which means one can be dominated by the other.

This fair splitting criterion explicitly considers survival time and censorship information now gives us a means to induce fair decision trees in the presence of censorship, and thus build Random Forests under uncertainty. We emphasize that this approach to promote fair splitting has no tunable parameters as given. This is to retain the desirable property of decision trees in that they often "just work" [50].

### C. Survival Universal Random Forests

With the tailored fair splitting criterion specifically accounts for censoring, we now introduce a corresponding learning

algorithm, *Survival Universal Random Forests (SURF)*, following the general idea of random forests to generate tailored forecasts while providing fair risk predictions. The pseudocode of SURF is illustrated in Algorithm 1.

---

**Algorithm 1** SURF Algorithm

---

**Input:** Censored and discriminated dataset $X$,
   The minimum unique events for splitting $d_0$,
   Ensemble size B
**Output:** Survival universal ensemble SURF
 1: **for** $i \in 1, \cdots, B$ **do**
 2:   $X^{(i)} \leftarrow$ A bootstrap sample from $X$
 3:   **while** Leaf $l$ with more than $d_0$ unique events **do**
 4:     $F^{(i)} \leftarrow$ A subset of the original features
 5:     $f \leftarrow$ The highest $S_\vee(f)$ for $f \in F^{(i)}$
          based on $X^{(i)}$
 6:     Split $l$ on $f$
 7:   **end while**
 8:   SURF= SURF $\cup$ $SURF_i$
 9: **end for**
10: **return** The learned SURF

---

The algorithm specifically works as follows: for each ensemble member of SURF, a bootstrap sample from the censored and discriminated dataset $X$ is first selected where $X^{(i)}$ denotes the $i$th bootstrap (line 1-2). When growing the ensemble member, the candidate splitting attributes at each node are restricted to some randomly selected subset of the whole feature space to decrease the correlation between trees in the ensemble (line 4). From fairness point of view, the randomness introduced by such restriction also limits the biased correlation inherited in the data. Line 5 then performs the splitting test related calculation based on the newly proposed USD to maximize the survival difference of all demographic groups rather than certain population. When the best splitting attribute in the subset has been selected, line 6 splits the node, causing each ensemble learner to grow. As selecting which feature to split is the most computationally expensive aspect of decision tree construction, narrowing the set of features further decreases the computational requirements of SURF. SURF is in full size when all terminal nodes of each base learner, *e.g*, $SURF_i$, contain less than $d_0 > 0$ unique event of interest (line 3).

With SURF being inducted, it is ready for fair risk predictions. SURF predicts the risk as the cumulative hazard function $H(t|x)$ to have a direct interpretation of the expected number of events, and it is the intermediate function between hazard and survival functions for direct derivation when needed. Formally, the risk score is estimated by the Nelson-Aalen estimator [52] as:

$$r(t|x) = \sum_{j \leq t} \frac{d_j}{n_j} \qquad (12)$$

where $d_j$ and $n_j$ represent the number of individuals experiencing the events and have not experienced the event at time

$j$ respectively, and $t$ is evaluated as the last event time. In response to cases within the same node sharing identical class label in non-censoring trees, all individuals within the node of SURF have the same risk score which is used for final risk predictions.

### D. Interplay between Fairness and Censorship

Armed with the previously introduced fairness notions and debiasing algorithm explicitly considering survival time and censoring information, quantifying and guaranteeing unfairness in the presence of censorship becomes feasible. In addition, it is also of practitioners and policy makers' interests to have an insight into the fairness-performance. However, due to censorship, existing fairness theoretical analyses are not applicable to understand the interplay between fairness and censorship. Here, we further provide a versatile geometric formalism to study fairness under uncertainty.

To this end, we formulate the fairness diagnostic as a fairness with uncertainty confusion tensor encoding the information needed to study the fair discriminative ability, *i.e.*, concordance imparity, of the model, which provides a general perspective for understanding unfairness under uncertainty. In specific, the fairness with uncertainty confusion tensor is the stack of confusion matrices for the sensitive attribute $G$, as shown in Table I. For simplicity, we assume that $G$ is a binary attribute, *i.e.*, $G \in \{0, 1\}$.

TABLE I. The fairness with uncertainty confusion tensor, showing the two planes corresponding to the confusion matrix for each of the favored (G = 1) and deprived groups (G = 0).

| G = 0 | r > r' | r < r' |
|---|---|---|
| t > t' & δ' = 0 | $N_0^a$ | $N_0^b$ |
| t > t' & δ' = 1 | $M_0^a$ | $M_0^b$ |
| t' > t & δ = 0 | $N_0^c$ | $N_0^d$ |
| t' > t & δ = 1 | $M_0^c$ | $M_0^d$ |

| G = 1 | r > r' | r < r' |
|---|---|---|
| t > t' & δ' = 0 | $N_1^a$ | $N_1^b$ |
| t > t' & δ' = 1 | $M_1^a$ | $M_1^b$ |
| t' > t & δ = 0 | $N_1^c$ | $N_1^d$ |
| t' > t & δ = 1 | $M_1^c$ | $M_1^d$ |

Let us denote the elements of fairness with uncertainty confusion tensor as $N_G^a$, $N_G^b$, $N_G^c$ and $N_G^d$ as well as $M_G^a$, $M_G^b$, $M_G^c$ and $M_G^d$ (abbreviations in align with Equation 6), each element with subscripts indicating the value of $G$. We further denote $P_G = M_G^a + M_G^b + M_G^c + M_G^d$, $I_G = N_G^a + N_G^b + N_G^c + N_G^d$ and $C_G = M_G^b + M_G^c$ be the number of pairwise individuals in each group $G \in \{0, 1\}$ pertaining permissible, impermissible and concordant pairs, respectively. Assume $P_G$, $I_G$ and $C_G$ are known constants. Unraveling the fairness with uncertainty confusion tensor, we can formulate concordance imparity as:

$$C_\triangle = \left| \frac{C_0}{P_0} - \frac{C_1}{P_1} \right| \tag{13}$$

We show below that how fair/unfair the model could appear in the presence of censorship, *i.e.*, how censorship could affect fairness notion calculation, based on "floor" outcome and "ceiling" outcome possible.

Let us first consider the *"floor" outcome*– all impermissible pairs, *e.g.*, $I_G$, become permissible with all censored individuals experience immediate event of interest, *i.e.*, $\delta$ becomes 1 for those individuals with $\delta = 0$ and their censored time $t$ now becomes the actual event time. Under this floor circumstance the CI becomes:

$$\lfloor C_\triangle \rfloor = \left| \frac{C_0 + I_0^{con}}{P_0 + I_0} - \frac{C_1 + I_1^{con}}{P_1 + I_1} \right| \tag{14}$$

where $I_0^{con} = N_0^b + N_0^c$ and $I_1^{con} = N_1^b + N_1^c$ representing the additional concordant pairs under this floor circumstance. Now we are ready to discuss different sub-scenarios possible affecting fairness notion calculation to provide practitioners and stakeholders more insights.

- *Sub-scenario 1:* $\frac{C_0}{P_0} = \frac{I_0^{con}}{I_0}$ *and* $\frac{C_1}{P_1} = \frac{I_1^{con}}{I_1}$.
  In this case, $\lfloor C_\triangle \rfloor = \left| \frac{C_0 + \frac{C_0 * I_0}{P_0}}{P_0 + I_0} - \frac{C_1 + \frac{C_1 * I_1}{P_1}}{P_1 + I_1} \right| = \left| \frac{C_0(1 + \frac{I_0}{P_0})}{P_0 + I_0} - \frac{C_1(1 + \frac{I_1}{P_1})}{P_1 + I_1} \right| = \left| \frac{C_0(1 + \frac{I_0}{P_0})}{P_0(1 + \frac{I_0}{P_0})} - \frac{C_1(1 + \frac{I_1}{P_1})}{P_1(1 + \frac{I_1}{P_1})} \right| = \left| \frac{C_0}{P_0} - \frac{C_1}{P_1} \right|$
  That is to say CI stays unchanged in this case. In practice, this suggests that the independence of model prediction with censorship. Practitioners and stakeholders can therefore interpret the level of unfairness of the model under consideration with more confidence.
- *Sub-scenario 2:* $\frac{C_0}{P_0} = 0$ *and* $\frac{C_1}{P_1} = 1$.
  In this case, the uncertainty affect the fairness calculation the most and could lead to the most deviated CI score from the original CI score in the presence of censorship. Sufficient care must be taken to ensure the extra unfairness possible as the result of the censorship when deploying social sensitive applications.
- *Sub-scenario 3:* $\frac{C_0}{P_0} = 0$ *and* $\frac{C_1}{P_1} = 0$.
  In contrast to the sub-scenario 2, the uncertainty, in this case, does affect the fairness calculation but uniform in direction. Practitioners and stakeholders can have this in mind considering real-world's socially sensitive impacts.
- *Sub-scenario 4: assume certain distribution with* $\frac{C_0}{P_0}$ *and* $\frac{C_1}{P_1}$, *for example the binomial distribution.*
  The assumed distribution on censorship in conjunction with existing CI score in the presence of censorship can therefore provide additional confidence interval with respect to the range of the CI score in the presence of censorship.

In contrast to the previous "floor" outcome, the second *"ceiling" outcome* also presents all impermissible pairs become permissible but all censored individuals now have an infinite event time of interest instead, *i.e.*, $\delta$ also becomes

1 for those individuals with $\delta = 0$ but their censored time $t$ becomes infinite instead. Under this floor circumstance the concordant elements of $\lceil C_\triangle \rceil$ change to the other way around, *i.e.*, $N_0^{con} = N_0^a + N_0^d$ and $N_1^{con} = N_1^a + N_1^d$ with all previously impermissible pairs, *e.g.*, $I_G$, become permissible as well. In addition, the censorship of longer survival time also makes difference and should be considered to form the new uncertainty tensor. Other than these, the previous discussion on "floor" outcome still applies and can be analyzed similarly to help the practical application deployment.

TABLE II. Summary of datasets used in experiments.

| Characteristics \ Datasets | ROSSI | COMPAS | KKBOX |
|---|---|---|---|
| Sample # | 432 | 10,325 | 2,814,735 |
| Censored # | 318 | 7,558 | 975,834 |
| Censored Rate | 0.736 | 0.732 | 0.347 |
| Feature # | 9 | 14 | 18 |
| Sensitive Attribute | race | race | gender |
| Sensitive Value | black | African American | female |

## V. EXPERIMENTS

### A. Datasets

We validate our proposed SURF on three real-world censored datasets with socially sensitive concerns and diverse characteristics[1]. Table II provides a summary description of them. Note that survival time and censoring information are explicitly included in our study to specifically account for censorship.

The **ROSSI** dataset pertains to persons convicted then released from Maryland state prisons, who were followed up for one year after release [53]. Among the total of 432 released convicts engaging in an experimental treatment, half of them were given financial aid and the rest of them did not receive aid. The task is then to predict the reoffending risk score of convicted criminals as described by 9 features, including the sensitive attribute "race" with black being the deprived group and other being the favored group.

The landmark algorithmic unfairness **COMPAS** dataset has a similar learning task but is considerably larger with 10,325 convicts from Broward County [54]. The COMPAS dataset consists of 14 features with "race" being the sensitive attribute and the African American defining the deprived group in contrast to other as the favored group. The COMPAS dataset also shares a similar 73% censored rate with the ROSSI dataset.

The **KKBOX** dataset is created from the WSDM-KKBOX's Churn Prediction Challenge 2017 [55]. Specially, the task is to predict a user's risk score of canceling his/her streaming music subscription from KKBOX. Compared with the previous two datasets, this is the largest number dataset we experiment with along with a lower censored rate of 34.7%. The total 2,814,735

---
[1]The code and data are publicly available at
https://github.com/vanbanTruong/censoredFairness

TABLE III. Evaluation results of different models with the best results marked in bold.

| Datasets | Metrics / Method | CI% | Fair Calibration | C-index% | Brier Score% | Time-dependent AUC% |
|---|---|---|---|---|---|---|
| ROSSI | IDCPH | 15.31 | Not fair calibrated | 52.28 | 18.73 | 77.32 |
| | GDCPH | 9.32 | **Fair calibrated** | 59.34 | 22.87 | 78.51 |
| | CPH | 11.43 | Not fair calibrated | 64.24 | 17.67 | 77.12 |
| | RSF | 16.53 | Not fair calibrated | 65.56 | 15.12 | 79.32 |
| | DeepSurv | 12.32 | Not fair calibrated | 66.67 | 14.71 | 77.17 |
| | SURF | **5.8** | **Fair calibrated** | **69.33** | **12.33** | **79.39** |
| COMPAS | IDCPH | 25.18 | Not fair calibrated | 62.16 | 25.03 | 63.78 |
| | GDCPH | 11.77 | **Fair calibrated** | 72.16 | 16.32 | 66.21 |
| | CPH | 22.43 | Not fair calibrated | 69.24 | 20.35 | 65.15 |
| | RSF | 25.32 | Not fair calibrated | 72.61 | 15.62 | 71.76 |
| | DeepSurv | 16.72 | Not fair calibrated | 75.12 | **13.42** | 71.83 |
| | SURF | **8.31** | **Fair calibrated** | **76.43** | 13.72 | **72.03** |
| KKBOX | IDCPH | 17.79 | Not fair calibrated | 72.61 | 21.23 | 69.73 |
| | GDCPH | 14.98 | **Fair calibrated** | 79.45 | 19.92 | 73.03 |
| | CPH | 18.91 | Not fair calibrated | 80.02 | 18.17 | 72.95 |
| | RSF | 21.14 | Not fair calibrated | 82.32 | 14.24 | 78.18 |
| | DeepSurv | 20.66 | Not fair calibrated | **83.01** | 14.33 | 80.71 |
| | SURF | **12.57** | **Fair calibrated** | 82.72 | **13.21** | **81.21** |

users are described using the same 18 features developed by the winning Kaggle team, and gender is used as the sensitive attribute with female being the sensitive value.

### B. Experimental Setup

We compare our method with the representative survival models and recent survival models with fairness constraints based on a set of typical survival analysis evaluation metrics as well as our proposed fairness notions considering censorship measuring their fair risk prediction capabilities.

**Evaluation Metrics.** In addition to the proposed fairness measures considering censorship (c.f. Section IV-A), in our evaluation we also report the following widely used survival analysis evaluation metrics: i) the *C-index*, which is an accuracy measure equals to the area under ROC curve (AUC) in the absence of censorship [56], ii) the *Brier score* in measuring the mean squared error between the probability estimations assigned to possible outcomes and the actual outcome [57], iii) the *Time-dependent AUC*, which is a function of time testing the discriminative capability of a model when distinguishing individuals experienced the event of interest from those have not experienced till time $t$ [58]. Among them, while the model with a lower Brier score is desired, the higher C-index and Time-dependent AUC scores, the better.

**Comparison Methods.** The proposed SURF is designed to address highly underexplored discrimination in the presence of censorship. To evaluate its design, we consider baselines from the following four perspectives: i) two recent proportional hazards assumption based fair survival models IDCPH and GDCPH [21], which are the only works for fair survival analysis problem to the best of our knowledge (note that only the most competitive ones are discussed among different variants proposed therein), ii) the mostly widely used CPH for survival analysis [38], iii) the state of the art random

forests on survival analysis RSF [43] and deep neural network based DeepSurv [39] as additional baselines. Other competing fairness methods are not considered as none of them is capable of addressing fairness in the presence of censoring.

### C. Results

We present both quantitative and qualitative experiment results that illustrate and confirm the utility of SURF in making fair decisions under uncertainty.

**Fair risk predictions.** A set of fairness metrics in the presence of censorship and survival analysis evaluations are tested on different models with the 5-fold cross validation results summarized in Table III.

As can be seen clearly our new SURF dominates all other baselines in terms of minimizing discrimination while maintaining a competitive predictive performance, which verifies the necessity of its debiasing design while accounting for censorship. On the other hand, the lack of consideration for survival time and censoring information as well as task-specific similarity metric among individuals result in the inferior performances of other baselines, although they are proposed for fairness in censoring settings. This also verifies our previous discussion that fairness in the presence of censorship cannot be trivially solved by a simple combination of existing techniques.

As p-values are not intended to be ranked, the fair calibration cannot rank models besides suggesting whether one model is calibrated across different demographic communities. We therefore further visualize the predicted probabilities against the true probability observed by KM to qualitatively evaluate the fairness of such predictions, as shown in Figure 1. Note that for space constraint, only the fair calibration plots for deprived community is presented. As we can see, the dark

Fig. 1. The binned fair calibration for deprived community, graphing the predicted probabilities against the true probability in deciles. The color of the bar represents different methods with the red bar representing the true probability observed by KM.
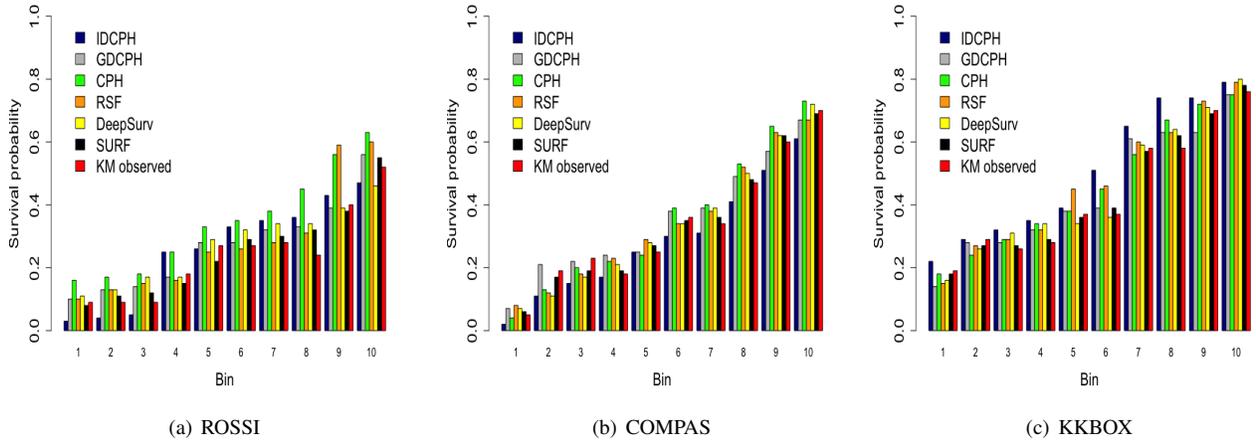
(a) ROSSI    (b) COMPAS    (c) KKBOX

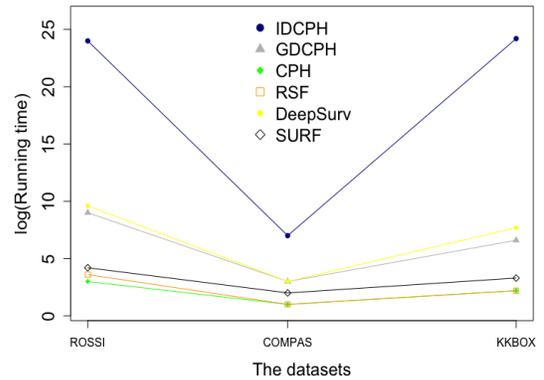TABLE IV. Predictive performance confusion matrix of SURF.

| Datasets | C-index% deprived | | C-index% favored | | Brier Score% deprived | | Brier Score% favored | | Time-dependent AUC% deprived | | Time-dependent AUC% favored | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SURF- | SURF | SURF- | SURF | SURF- | SURF | SURF- | SURF | SURF- | SURF | SURF- | SURF |
| ROSSI | 51.71 | 65.32 | 74.24 | 71.12 | 21.03 | 16.23 | 9.87 | 10.17 | 69.98 | 73.65 | 84.87 | 82.17 |
| COMPAS | 54.82 | 70.22 | 80.14 | 78.53 | 18.76 | 16.12 | 7.66 | 11.65 | 62.81 | 65.17 | 77.62 | 75.24 |
| KKBOX | 64.53 | 72.89 | 85.67 | 85.46 | 18.87 | 15.17 | 7.16 | 9.12 | 72.31 | 77.45 | 85.87 | 84.66 |

bars, are roughly the same height as the red ones, suggesting SURF's predicted probabilities are representative of the deprived community's true probabilities.

In addition, we also perform an ablation study, testing whether the debiasing component of SURF worked to benefit the improved performance of the deprived group. To this end, we further break down the predictive performance by the sensitive value that defines the deprived community and favored community and with and without our fairness constraints (*e.g.*, SURF and SURF-). The predictive performance confusion matrix are shown in Table IV. As one can see, SURF attends to the deprived community with improved predictive ability and characterization, which verifies its anti-discrimination capability. In addition, the improved overall predictive performance also shows the merit of such anti-discrimination design for fair risk predictions as a whole.

**Time Complexity.** We illustrate and analyze the time complexity of our proposed SURF against other baseline methods, as shown in Figure 2, where the y-axis is the logarithm of the runtime in seconds for ROSSI dataset and minutes for the other two datasets. As one can see, SURF enjoys the merit of computational efficiency as it follows the idea of tree based method to decrease the computational requirements in practice, while the other baseline IDCPH becomes computational expensive, thus limiting its applicability.

Fig. 2. The time complexity comparison of all the methods.



## VI. CONCLUSION

This work is motivated by the increasing attention on the issue of discriminatory AI behaviors, with the aim to provide fair and accuracy predictions in the presence of censorship. Different from existing fair survival methods, we propose two censored-specific fairness notions and a new debiasing algorithm to specifically account for fairness under censorship. The positive results of conducted experiments show the anti-discrimination capability of our proposed method. The proposed technique is expected to be versatile in alleviating bias

in various socially sensitive applications (*e.g.*, the allocations of health resources, personalized marketing and recidivism prediction instrument).

## REFERENCES

[1] A. Beutel, J. Chen, T. Doshi, H. Qian, A. Woodruff, C. Luu, P. Kreitmann, J. Bischof, and E. H. Chi, "Putting fairness principles into practice: Challenges, metrics, and improvements," *AIES*, 2019.

[2] D. Meyer, "Amazon reportedly killed an ai recruitment system because it couldn't stop the tool from discriminating against women. fortune, october 10," 2018.

[3] M. Skirpan and M. Gorelick, "The authority of "fair" in machine learning," 2017.

[4] K. Grace, J. Salvatier, A. Dafoe, B. Zhang, and O. Evans, "When will ai exceed human performance? evidence from ai experts," *Journal of Artificial Intelligence Research*, vol. 62, pp. 729–754, 2018.

[5] S. Barocas, M. Hardt, and A. Narayanan, "Fairness in machine learning," *Nips tutorial*, vol. 1, p. 2, 2017.

[6] R. Binns, "Fairness in machine learning: Lessons from political philosophy," in *Conference on Fairness, Accountability and Transparency*, 2018, pp. 149–159.

[7] V. Chang, "An ethical framework for big data and smart cities," *Technological Forecasting and Social Change*, vol. 165, p. 120559, 2021.

[8] C. Chen and R. Wong, "Black patients miss out on promising cancer drugs—propublica. 2018," 2019.

[9] S. Barocas and A. D. Selbst, "Big data's disparate impact," *California Law Review*, vol. 104, no. 3, p. 671, 2016.

[10] A. Agarwal, A. Beygelzimer, M. Dudík, J. Langford, and H. Wallach, "A reductions approach to fair classification," in *International Conference on Machine Learning*. PMLR, 2018, pp. 60–69.

[11] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel, "Fairness through awareness," in *Proceedings of the 3rd innovations in theoretical computer science conference*, 2012, pp. 214–226.

[12] W. Zhang and J. Weiss, "Fair decision-making under uncertainty," in *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2021.

[13] W. Zhang and J. C. Weiss, "Longitudinal fairness with censorship," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 11, 2022, pp. 12 235–12 243.

[14] T. G. Clark, M. J. Bradburn, S. B. Love, and D. G. Altman, "Survival analysis part i: basic concepts and first analyses," *British journal of cancer*, vol. 89, no. 2, pp. 232–238, 2003.

[15] A. Beutel, J. Chen, Z. Zhao, and E. H. Chi, "Data decisions and theoretical implications when adversarially learning fair representations," *arXiv preprint arXiv:1707.00075*, 2017.

[16] S. Hajian, F. Bonchi, and C. Castillo, "Algorithmic bias: From discrimination discovery to fairness-aware data mining," in *Proceedings of the SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 2125–2126.

[17] W. Zhang and J. Weiss, "Fairness with censorship and group constraints," *Knowledge and Information Systems*, 2022.

[18] W. Zhang, J. C. Weiss, S. Zhou, and T. Walsh, "Fairness amidst non-iid graph data: A literature review," *arXiv preprint arXiv:2202.07170*, 2022.

[19] F. Kamiran and T. Calders, "Data preprocessing techniques for classification without discrimination," *Knowledge and Information Systems*, vol. 33, no. 1, pp. 1–33, 2012.

[20] M. B. Zafar, I. Valera, M. Gomez Rodriguez, and K. P. Gummadi, "Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment," in *Proceedings of the 26th international conference on world wide web*, 2017, pp. 1171–1180.

[21] K. N. Keya, S. Pan, I. Stockwell, and J. Foulds, "Equitable allocation of healthcare resources with fair cox models," in *AAAI Fall Symposium on AI in Government and Public Sector*, 2020.

[22] M. J. Bradburn, T. G. Clark, S. B. Love, and D. G. Altman, "Survival analysis part ii: multivariate data analysis–an introduction to concepts and methods," *British journal of cancer*, vol. 89, no. 3, pp. 431–436, 2003.

[23] Y. Bechavod, C. Jung, and S. Z. Wu, "Metric-free individual fairness in online learning," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

[24] F. Kamiran and T. Calders, "Classifying without discriminating," in *2nd International Conference on Computer, Control and Communication*, 2009, pp. 1–6.

[25] I. Žliobaite, F. Kamiran, and T. Calders, "Handling conditional discrimination," in *2011 IEEE 11th International Conference on Data Mining*. IEEE, 2011, pp. 992–1001.

[26] T. Calders, F. Kamiran, and M. Pechenizkiy, "Building classifiers with independency constraints," in *ICDMW*, 2009, pp. 13–18.

[27] W. Zhang and A. Bifet, "Feat: A fairness-enhancing and concept-adapting decision tree classifier," in *International Conference on Discovery Science*. Springer, 2020, pp. 175–189.

[28] W. Zhang *et al.*, "Flexible and adaptive fairness-aware learning in non-stationary data streams," in *IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI)*, 2020, pp. 399–406.

[29] W. Zhang and L. Zhao, "Online decision trees with fairness," *arXiv preprint arXiv:2010.08146*, 2020.

[30] M. Hardt, E. Price, N. Srebro *et al.*, "Equality of opportunity in supervised learning," in *Advances in neural information processing systems*, 2016, pp. 3315–3323.

[31] B. Fish, J. Kun, and Á. D. Lelkes, "A confidence-based approach for balancing fairness and accuracy," in *SDM*, 2016, pp. 144–152.

[32] J. Zeng, B. Ustun, and C. Rudin, "Interpretable classification models for recidivism prediction," *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, vol. 180, no. 3, pp. 689–722, 2017.

[33] Z. Liu, R. Wang, N. Japkowicz, D. Tang, W. Zhang, and J. Zhao, "Research on unsupervised feature learning for android malware detection based on restricted boltzmann machines," *Future Generation Computer Systems*, vol. 120, pp. 91–108, 2021.

[34] W. Zhang, J. Wang, D. Jin, L. Oreopoulos, and Z. Zhang, "A deterministic self-organizing map approach and its application on satellite data based cloud type classification," in *IEEE International Conference on Big Data (Big Data)*, 2018.

[35] W. Zhang, J. Tang, and N. Wang, "Using the machine learning approach to predict patient survival from high-dimensional survival data," in *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2016.

[36] W. Zhang and J. Wang, "Content-bootstrapped collaborative filtering for medical article recommendations," in *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2018.

[37] W. Zhang, L. Zhang, D. Pfoser, and L. Zhao, "Disentangled dynamic graph deep generation," in *Proceedings of the SIAM International Conference on Data Mining (SDM)*, 2021, pp. 738–746.

[38] D. R. Cox, "Regression models and life-tables," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 34, no. 2, pp. 187–202, 1972.

[39] J. L. Katzman, U. Shaham, A. Cloninger, J. Bates, T. Jiang, and Y. Kluger, "Deepsurv: personalized treatment recommender system using a cox proportional hazards deep neural network," *BMC medical research methodology*, vol. 18, no. 1, pp. 1–12, 2018.

[40] X. Tang, X. Huang *et al.*, "Cognitive visual commonsense reasoning using dynamic working memory," in *International Conference on Big Data Analytics and Knowledge Discovery*. Springer, 2021.

[41] I. Bou-Hamad, D. Larocque, H. Ben-Ameur *et al.*, "A review of survival trees," *Statistics surveys*, vol. 5, pp. 44–71, 2011.

[42] X. Tang, L. Zhang *et al.*, "Using machine learning to automate mammogram images analysis," in *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2020, pp. 757–764.

[43] H. Ishwaran, U. B. Kogalur, E. H. Blackstone, M. S. Lauer *et al.*, "Random survival forests," *Annals of Applied Statistics*, vol. 2, no. 3, pp. 841–860, 2008.

[44] P. Wang, Y. Li, and C. K. Reddy, "Machine learning for survival analysis: A survey," *ACM Computing Surveys (CSUR)*, vol. 51, no. 6, pp. 1–36, 2019.

[45] R. G. Miller Jr, *Survival analysis*. John Wiley &amp; Sons, 2011, vol. 66.

[46] A. Latouche, A. Allignol, J. Beyersmann, M. Labopin, and J. P. Fine, "A competing risks analysis should report results on all cause-specific hazards and cumulative incidence functions," *Journal of clinical epidemiology*, vol. 66, no. 6, pp. 648–653, 2013.

[47] D. W. Hosmer and S. Lemesbow, "Goodness of fit tests for the multiple logistic regression model," *Communications in statistics-Theory and Methods*, vol. 9, no. 10, pp. 1043–1069, 1980.

[48] R. B. D'Agostino and B.-H. Nam, "Evaluation of the performance of survival analysis models: discrimination and calibration measures," *Handbook of statistics*, vol. 23, pp. 1–25, 2003.

[49] J. Ranstam and J. Cook, "Kaplan–meier curve," *British Journal of Surgery*, vol. 104, no. 4, pp. 442–442, 2017.

[50] J. Han, M. Kamber, and J. Pei, "Data mining: concepts and techniques, waltham, ma," 2012.

[51] J. M. Bland and D. G. Altman, "The logrank test," *Bmj*, vol. 328, no. 7447, p. 1073, 2004.

[52] Ø. Borgan, "Nelson–aalen estimator," *Wiley StatsRef: Statistics Reference Online*, 2014.

[53] J. Fox, M. S. Carvalho *et al.*, "The rcmdrplugin. survival package: Extending the r commander interface to survival analysis," *Journal of Statistical Software*, vol. 49, no. 7, pp. 1–32, 2012.

[54] J. Angwin, J. Larson, S. Mattu, and L. Kirchner, "There's software used across the country to predict future criminals," *And it's biased against blacks. ProPublica*, 2016.

[55] H. Kvamme, Ø. Borgan, and I. Scheel, "Time-to-event prediction with neural networks and cox regression," *Journal of Machine Learning Research*, vol. 20, no. 129, pp. 1–30, 2019.

[56] F. E. Harrell, R. M. Califf, D. B. Pryor, K. L. Lee, and R. A. Rosati, "Evaluating the yield of medical tests," *Jama*, vol. 247, no. 18, pp. 2543–2546, 1982.

[57] G. W. Brier and R. A. Allen, "Verification of weather forecasts," in *Compendium of meteorology*. Springer, 1951, pp. 841–848.

[58] L. E. Chambless and G. Diao, "Estimation of time-dependent area under the roc curve for long-term risk prediction," *Statistics in medicine*, vol. 25, no. 20, pp. 3474–3486, 2006.