# Equalised Odds Is Not Equal Individual Odds:
# Post-processing for Group and Individual Fairness

EDWARD SMALL, RMIT University, Australia

KACPER SOKOL, ETH Zurich, Switzerland

DANIEL MANNING, RMIT University, Australia

FLORA D. SALIM, UNSW Sydney, Australia

JEFFREY CHAN, RMIT University, Australia

Group fairness is achieved by equalising prediction distributions between protected sub-populations; individual fairness requires treating similar individuals alike. These two objectives, however, are incompatible when a scoring model is calibrated through *discontinuous* probability functions, where individuals can be randomly assigned an outcome determined by a fixed probability. This procedure may provide two similar individuals from the same protected group with classification odds that are disparately different – a clear violation of individual fairness. Assigning unique odds to each protected sub-population may also prevent members of one sub-population from ever receiving the chances of a positive outcome available to individuals from another sub-population, which we argue is another type of unfairness called *individual odds*. We reconcile all this by constructing continuous probability functions between group thresholds that are constrained by their Lipschitz constant. Our solution preserves the model's predictive power, individual fairness and robustness while ensuring group fairness.

## 1 INTRODUCTION

Predictive models that output a score or probability for a multi-dimensional input, i.e., scoring functions, are a common tool in automated decision-making [12, 13]. Binary classification is a popular realisation of this paradigm, where a threshold is placed on a score to produce a decision; among others, it can be found in school examinations where individual answers are condensed into a grade that translates to a pass/fail mark [37], or banking where the history of personal finances is compressed into a credit score that captures one's likelihood of defaulting on a loan [27]. Many

Authors' Contact Information: Edward Small, RMIT University, Melbourne, Australia, edward.small@student.rmit.edu.au; Kacper Sokol, ETH Zurich, Zurich, Switzerland, kacper.sokol@inf.ethz.ch; Daniel Manning, RMIT University, Melbourne, Australia, daniel.manning@student.rmit.edu.au; Flora D. Salim, UNSW Sydney, Sydney, Australia, flora.salim@unsw.edu.au; Jeffrey Chan, RMIT University, Melbourne, Australia, jeffrey.chan@rmit.edu.au.
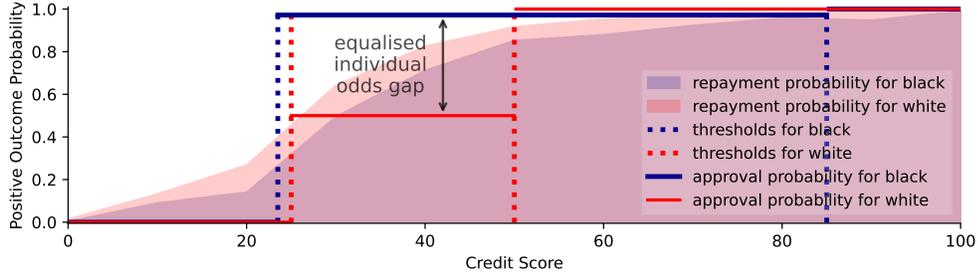
Fig. 1. Two-threshold fixed randomisation [22] applied to probabilities (y-axis) output by a loan repayment classifier built upon credit scores (x-axis). It satisfies equalised odds for the binary protected attribute *race* (black and white) by using it to assign approval probabilities, but results in discontinuities that violate individual fairness and create a gap between group-specific *individual odds*.

such applications, especially in high stakes domains like healthcare, finance and judiciary, are coming under increased scrutiny given their potential harm to society – predictive models deployed in these contexts are expected to be accurate, robust, fair and explainable. These four desiderata, however, are often at odds. Improving utility, i.e., predictive power, of a model may entail increasing its complexity at the expense of interpretability and robustness, e.g., due to overfitting [4, 46]. Similarly, equalising errors between protected groups to ensure fairness may require sacrificing utility and impairing other notions of fairness [14, 41].

In this paper we focus on the latter scenario, where (protected) sub-populations are treated differently, thus unfairly, due to persistent historical biases [10], training data under-representation [11] and greedy optimisation of an objective function. Correcting for these biases is often challenging as it requires detailed knowledge of the data domain and the input space. One popular solution to this problem, which we study here, is threshold optimisation under fairness constraints when dealing with multiple protected groups. This method relies on calculating unique decision functions based on scores for each protected group in order to satisfy a given fairness constraint, e.g., demographic parity [39].

We re-examine this approach, as finding a set of thresholds for a given score function that satisfy multiple fairness constraints – such as equalised odds [3] – is often impossible if only using a collection of single thresholds. Instead, a decision function that is optimal with respect to a definition of group fairness selected by the model owner is derived directly from the scoring function using a pair of thresholds for each group. Outputs that fall between the thresholds are allocated a random decision based on a fixed probability parameter – a procedure called *fixed randomisation* – which, while effective, exhibits a number of shortcomings demonstrated by Figure 1 and discussed later in Section 3. Using a fixed randomisation parameter is suboptimal for both the entities that create the model (owners) and those whose case is being decided by the model (users) because:

(1) if the scoring function is accurate, the decision function cannot leverage this in the intervals between the thresholds;

(2) even if the scoring function is individually fair, the step-based decision function is not, e.g., users whose scores are just under a threshold are treated very differently to those who are barely above it, despite their scores being similar (commonly referred to as the *threshold effect* [32]); and

(3) users from one protected group may be unable to access the odds of positive classification offered to another group (equalised *individual odds* unfairness).

Consider the fixed randomisation solution shown in Figure 1, which satisfies equalised odds for the two values of the protected attribute *race* in a loan allocation setting using fixed randomisation. For example, a *white* user with a credit

score of 49.5 is assigned the same odds (50%) of receiving a loan as a *white* user whose credit score is 25 despite the latter being 6.6 times more likely to default than the former. This stands in stark contrast to a *white* user with a credit score of 24.5 (just below the threshold) who has no chance of getting a loan despite being only 1.03 times more likely to default than the aforementioned *white* user with a credit score of 25. This case study illustrates that an increase in credit score – and therefore an increase in the likelihood of repaying a loan – is not reflected in the final decision for all scores except at the thresholds. In addition, while some *white* users have a 50% chance of receiving a loan and some *black* users have a 97.2% chance, these success odds are never offered to the other group; therefore, a *white* user will never be given a 97.2% chance of receiving a loan and vice versa. This disparity motivates a new notion of fairness – called *equalised individual odds* – which we outline in Definition 3.2 in Section 3.

We address these shortcomings by deriving a set of closed-form, continuous, monotonic functions (shown later in Figure 4) parameterised only by the thresholds and a probability parameter, making them easy to compute (Section 4.2). We show that these functions are constrained via a maximum derivative, preventing a change in score leading to a large shift in classification odds and thus maintaining individual fairness and softening the threshold effect (Section 4.3). Our approach enables the model owners to prioritise users with higher scores, better honouring the underlying score distribution as well as improving the transparency of the process. These properties incentivises users to increase their score as such an action improves their odds of a positive outcome – see Figure 3 for a direct comparison to Figure 1. We analyse our method in two case studies – through the lens of credit scoring for loan allocation in Section 5.1, and risk of recidivism in Section 5.2. For the credit scoring case study we seek equalised odds across the four values of the *race* attribute – non-Hispanic white (white), black, Hispanic and Asian – found in the *2003 TransUnion TransRisks Scores* (CreditRisk) data set [38], whereas for the recidivism case study we enforce equalised odds across a combination of two *races* – Caucasian and African America – and two *sexes* – male and female – found in the *2016 ProPublica Recidivism Risk Score* (COMPAS) data set [38]. In both cases, we show that individual fairness is improved while group fairness and accuracy are preserved. In summary, our contribution is threefold: (1) we demonstrate that fixed randomisation for group fairness violates individual fairness; (2) we derive a set of closed-form, continuous and monotonic probability functions; and (3) we show that these continuous curves preserve group fairness and improve performance while adhering to the constraint imposed by individual fairness.

## 2 PRELIMINARIES

### 2.1 Notation

We assume that the scalar scoring function $g : \mathcal{X} \mapsto \mathcal{R}$ takes individual instances and outputs a score $\mathcal{R} \subseteq \mathbb{R}$; $h : \mathcal{R} \mapsto \mathcal{Y}$, where $\mathcal{Y} \equiv \{0, 1\}$, is an arbitrary, possibly stochastic, binary decision function on $\mathcal{R}$ that maps the scores $R$ to predicted classes $\widehat{Y}$ according to a predetermined probability distribution $\mathbb{P}\{\widehat{Y} = 1 | R = r\}$. Lower case letters denote an individual instance from a sample, e.g., $\mathbf{x}$ is an instance in $X$. Functions denoted by Greek letters, such as $\zeta : \mathcal{R} \mapsto \mathcal{I}$ where $\mathcal{I} \equiv [0, 1]$, parameterise this probability based on scores, e.g., according to the Bernoulli distribution $h(r) \sim B(1, \zeta(r))$. Effectively, $h(r)$ is the probability that $\widehat{Y} = 1$ for $R = r$. Alternatively, for deterministic behaviour $\zeta$ can be defined by a single threshold $t \in \mathcal{R}$, where a score $r \geq t$ yields $h(r) = 1$ and $r < t$ yields $h(r) = 0$; this behaviour can be captured by the indicator function

$$\zeta(r) = \mathbb{1}_t(r) = \begin{cases} 0 & \text{if } r < t \\ 1 & \text{if } r \geq t \end{cases}.$$

One common realisation of this thresholding function is a binary probabilistic classifier, where $\mathcal{R} \equiv \mathcal{I}$ and $t = 0.5$. We therefore define the final *decision function* $f_h : \mathcal{X} \mapsto \mathcal{Y}$ as the composition $f_h = h \circ g$, where the subscript on $f$ indicates the composition function $h$ on the scoring function $g$. Additionally, capital letters refer to samples from spaces, such that $X$ is a sample from the space $\mathcal{X}$; $g(X) = R$ are the corresponding scores calculated by $g$ for the sample $X$; $f(X) = \widehat{Y}$ are classes predicted for all instances in the sample $X$; and $Y$ captures their ground truth labels. We denote the protected attribute as $A$, and consider the joint distribution $(R, A, Y)$. We make no assumptions on the type or shape of $\mathcal{X}$, nor on the construction of $g$ (the behaviour of which is discussed in Section 3).

## 2.2 Distance and Similarity Measures

Defining "similar individuals" can be challenging and is deeply rooted in the landscape and shape of the input space, the complexity of the problem, and the density and distribution of the training data $X$ within the space. Distances on metric spaces, regardless of their definition, must follow a set of axioms (outlined in Appendix A). This problem is also not strictly mathematical and depends highly on the context. Additionally, discrete or categorical data can be difficult to quantify and compare; for example, in a feature space of size $N$, how different is an unmarried individual from a married person, all other things being equal? One could argue that its importance depends on the size of $N$ – a large value of $N$ can dilute the importance of each individual feature. If we are trying to predict whether an individual has any children, however, this feature is of high importance regardless of the size of $N$. To best capture such dependencies, we can employ similarity graphs or bespoke distance metrics chosen based on the problem definition and the data set at hand.

Using tailor-made definitions of similarity, nonetheless, poses two issues: (1) it makes it difficult to compare results between experiments; and (2) the results are subject to the quality of the metric and its suitability for the problem at hand. We operate under the assumption that model inputs are inaccessible (simulating scenarios where data and model parameters are protected and/or private), thus we are only given scores, values of the protected attribute and the label (ground truth). For our work we therefore rely on generic distance metrics such as Euclidean, Hamming and Gower's distances. Note that we assume that changing the protected class $A$ for an individual is too large of a change to label the two instances as similar since this alteration entails using a different set of thresholds and probabilities in the final decision function. Examples of classical distance functions are presented in Appendix A.

## 2.3 Related Work

Group and individual fairness are two commonly considered categories [7]. **Group fairness** focuses on the statistical difference in outcomes between sub-populations determined by the values of a protected attribute $A$ [5]. The type of statistical outcome that a model owner may want to focus on is domain-specific, but measures closer to 0 are more desirable as this indicates no statistical difference between two groups. For a simple case of a binary protected feature $A = \{a, a'\}$ where $a \cap a' = \emptyset$, we can further differentiate two types of group fairness [9]:

**Outcome** Predictions are equalised in a set way across groups, e.g., demographic parity [1]:

$$\left| \mathbb{P}\{\widehat{Y} = 1 | A = a\} - \mathbb{P}\{\widehat{Y} = 1 | A = a'\} \right| = 0 .$$

An example of demographic parity may be in school admissions [44], where the distribution of admitted students should represent the distribution of the applicants for each value of $A$ (i.e., if applicants are 50% male and 50% female, admissions should reflect this pattern).

**Error Distribution** (In)correct classifications should be equalised in a predetermined way, e.g., using false negative
   rate:

$$\left| \mathbb{P}\{\widehat{Y} = 0 | A = a, Y = 1\} - \mathbb{P}\{\widehat{Y} = 0 | A = a', Y = 1\} \right| = 0 \ .$$

   An example of equalising false negative rate may be in the medical field, where false negatives could have dire
   consequences for a patient. Erring on the side of caution equally for all groups is therefore more preferable, up
   to a certain cost [31].

There are many ways in which group fairness can be operationalised, with different tasks and domains requiring a
specific constraint or a mixture thereof. In this paper, we mainly consider one of the strongest fairness constraints
called *equalised odds* [33], which is outlined in Definition 2.1.

DEFINITION 2.1 (EQUALISED ODDS). *A decision function $f : \mathcal{X} \mapsto \mathcal{Y}$ satisfies equalised odds with respect to a protected
attribute A if false positives and true positives are independent of the protected attribute:*

$$\left| \mathbb{P}\{\widehat{Y} = 1 | A = a, Y = y\} - \mathbb{P}\{\widehat{Y} = 1 | A = a', Y = y\} \right| = 0 \ \ \forall y \in \mathcal{Y} \ \ \forall a, a' \in A \ \ a \neq a' \ .$$

A large portion of fairness research in machine learning therefore focuses on equalising outcomes and errors between
users who belong to different protected groups, such as race or sex [28]. There are three distinct areas where fairness
can be injected into a data modelling pipeline:

**pre-processing** transforms the underlying training data such that signals and cross-correlations causing bias and
   discrimination are weakened [8];
**in-processing** incorporates fairness constraints directly into the optimisation objective [40]; and
**post-processing** alters the output of a decision-making process to mitigate bias of the underlying (fixed) model [25].

A variety of methods is needed as even when the scoring function is trained as "unaware" [15], and as such has no
knowledge of the value of the protected class $A$, $f$ can still become unfair. For example, the ground truth $Y$ may
be correlated with $A$ due to historical biases, some features in $X$ may act as a proxy, or the distribution/behaviour
of some features in $X$ may be different between sub-populations, causing a predictive model to under-perform for
under-represented groups. A different strand of work looks into fair data collection [43] and feature selection [20] as well
as fair learning procedures, e.g., adversarial learning [45]. In this paper, we focus on a popular class of post-processing
methods known as *threshold optimisation*. Our work builds directly upon the foundational method introduced by Hardt
et al. [22] by expanding and improving it along multiple dimensions.

   A slightly more nuanced view on fairness is the notion of "treating similar individuals similarly", known as **individual
fairness** [30]. In short, we look to impose a constraint on the distance between any two points (individuals) in the
input space against their distance in the output space [15]. We measure distance or similarity using distance functions
$d$ on the input ($d_{\mathcal{X}}$) and output ($d_{\mathcal{R}}$) spaces:

$$d_{\mathcal{R}} \left( g(\mathbf{x}_1), g(\mathbf{x}_2) \right) \leq L_{\mathcal{X}} d_{\mathcal{X}}(\mathbf{x}_1, \mathbf{x}_2) \implies \frac{d_{\mathcal{R}}(r_1, r_2)}{d_{\mathcal{X}}(\mathbf{x}_1, \mathbf{x}_2)} \leq L_{\mathcal{X}} \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X} \ , \tag{1}$$

where $r_k = g(\mathbf{x}_k)$ and $L_{\mathcal{X}} \geq 0$ is a *Lipschitz constant* (refer to Section 2.2 for a discussion of distance metrics). The
Lipschitz constant describes the maximal difference of the distance between two values in the input space with their
corresponding distance in the output space. Limiting $L_{\mathcal{X}}$ is usually done with a smoothing process, e.g., manifold
regularisation [6], or by constraining the optimisation of $g$ subject to a condition on the size of $L_{\mathcal{X}}$. The concept is
to assume that individuals with similar features (small $d_{\mathcal{X}}$) should appear close together in the output space (small
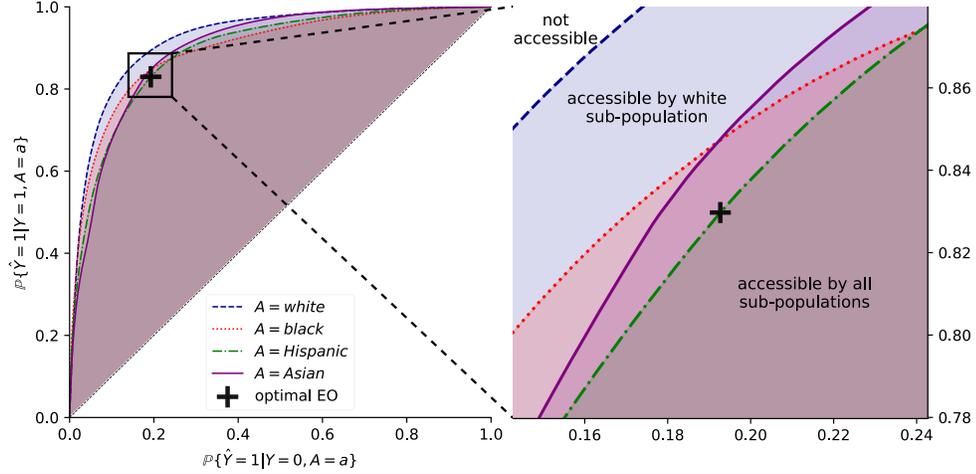
Fig. 2. ROC curves for the CreditRisk data set. The solution space for each (protected) group is given by all the points on their respective ROC curve when a single threshold is used. If we rely on multiple thresholds and randomisation, however, we expand the solution space to all the points on and below an ROC curve – represented for each group as the coloured area. A fair solution, according to equalised odds, is any set of thresholds and probabilities such that each group achieves equal true and false positives.

$d_{\mathcal{R}}$). Therefore, limiting the rate at which $g$ can change (i.e., its differential) in densely populated areas of the feature space can force $g$ to be smoother, hence more fair. $d_{\mathcal{X}}$ can be as simple as Gower's distance for mixed categorical and numerical features (see Appendix A), but ideally should be chosen appropriately for the problem at hand; since $\mathcal{R}$ is scalar, we can define $d_{\mathcal{R}}\left(g(\mathbf{x}_1), g(\mathbf{x}_2)\right) = |r_1 - r_2|$.

## 3  POST-PROCESSING FOR FAIRNESS WITH FIXED RANDOMISATION

In general, we expect that the higher the score $r$ output by a scoring function $g$ used for a predictive task, the "better" the outcome. This property is known as *positive orientation*, with *negative orientation* describing the opposite behaviour [19]. For example, if $\mathbf{x}_1$ and $\mathbf{x}_2$ are randomly drawn instances from $X$ used to calculate credit scores, and $\mathbf{x}_1$ has a higher credit score than $\mathbf{x}_2$ – i.e., $g(\mathbf{x}_1) > g(\mathbf{x}_2)$ – this relation implies that $\mathbf{x}_1$ is more likely to have healthier spending habits, thus making this person more likely to repay a loan (see Figure 8 in Appendix E). This does not need to be strictly true for all values, but should hold in general. In other words, we expect the receiver operating characteristic (ROC) curve – which expresses the (*false positive*, *true positive*) rates at different thresholds on $R$ – to at least be above the diagonal line from $(0, 0)$ to $(1, 1)$ and monotonic, i.e., never decreasing [18]. An ROC curve that is a straight line from $(0, 0)$ to $(1, 1)$ denotes a scoring function that is completely independent of $Y$, i.e., a trivial scoring function.

Optimising a scoring function $g$ with respect to complex definitions of fairness (such as equalised odds given by Definition 2.1) for multiple protected groups is more challenging than optimising $g$ for less strict fairness notions (e.g., demographic parity) due to the need to satisfy multiple constraints simultaneously. With post-processing, we assume that $g$ is *fixed* and *inaccessible*, i.e., a black box. We cannot therefore know or alter how scores are calculated from the input space, nor do we have access to the input space. This may be due to trade secrets or privacy concerns [24], and applies to credit scoring [23] among other domains. To achieve the desired notion(s) of fairness we therefore need to build an unbiased decision function $f_h$ upon $g$ by finding optimal thresholds for $h$ using *only* the joint distributions of $R$, $A$ and $Y$.

When cardinality of the protected attribute $A$ is 2, optimal equalised odds can be achieved by fixing a single threshold at any point where the ROC curves are equal. If there are multiple points where the curves meet, the optimal solution (lowest false positive and highest true positive rates) is the intersection closest to $(0, 1)$, i.e., the perfect model. Figure 2 shows the ROC curves stratified by a protected attribute $A$ (*race*) for a loan repayment prediction based on credit scores $R$ from the CreditRisk data set.

The challenge arises when the ROC curves do not touch or if $|A| > 2$. If the curves do not touch in $(0, 1) \times (0, 1)$, we can only satisfy equalised odds with a single threshold at the trivial points $(0, 0)$ or $(1, 1)$, i.e., assign the same outcome to all the scores. For $|A| > 2$ – e.g., where $A = A_1 \times A_2 \times \cdots \times A_n$ may be a Cartesian product of $n$ protected characteristics – it is highly unlikely for all the ROC curves to intersect at the same point (except for the trivial points). When using a single threshold, each group can only access false and true positive values that are **on their respective ROC curve** (shown in Figure 2 as the coloured curved lines). Using multiple thresholds and randomisation, however, allows each group to access all the points that are **below their respective ROC curve** and above the trivial scoring function (shown in Figure 2 as the coloured regions). The optimal point for equalised odds therefore becomes the point under all ROC curves that is closest to $(0, 1)$.

Hardt et al. [22] achieve equalised odds by setting group-specific thresholds $t_{y,a}$ – where $t_{0,a} \leq t_{1,a}$, so $y \in \{0, 1\}$ – that are applied to the scoring function $g$. If a score falls between the thresholds designated for the protected group $a$, it is assigned a class at random with a probability given by the parameter $p_a \in \mathcal{I}$. Since thresholds are group-specific, we define a threshold-based classification function $h_a : \mathcal{R} \mapsto \mathcal{Y}$, where the probability of $h_a(r) = 1$ is given by

$$\zeta_a(r) = p_a \mathbb{1}_{t_{0,a}}(r) + (1 - p_a)\mathbb{1}_{t_{1,a}}(r) , \tag{2}$$

for each protected sub-population $a$. In other words, $h_a(r) \sim B\big(1, \zeta_a(r)\big)$. We therefore define the final decision function as $f_{h_a} = h_a \circ g$, and Equation 2 gives us

$$\mathbb{P}\{f_{h_a}(\mathbf{x}) = 1 | A = a, X = \mathbf{x}\} = \begin{cases} 0 & \text{if } g(\mathbf{x}) < t_{0,a} \\ p_a & \text{if } g(\mathbf{x}) \in [t_{0,a}, t_{1,a}) \\ 1 & \text{if } g(\mathbf{x}) \geq t_{1,a} \end{cases} .$$

We call this *fixed randomisation*, as $r \in [t_{0,a}, t_{1,a})$ yields probability $p_a$ of $\widehat{Y} = 1$. Setting $p_a = 0$, $p_a = 1$ or $t_{0,a} = t_{1,a}$ is synonymous to using a single threshold. A visual example of fixed randomisation is provided in Figure 1.

Fixed randomisation is an effective approach to build a classifier $f_{h_a}$ based on a scoring function $g$ that satisfies group fairness such as equalised odds. This strategy, however, exhibits a number of undesired properties; most notably:

(i) it does not follow the general behaviour expected of a scoring function since all users who are subject to randomisation receive the same classification odds, no matter their score, but users whose scores are similar and near the thresholds are treated differently (ergo the example given in the introduction and shown in Figure 1);

(ii) even if $g$ is individually fair with well-defined $L_\mathcal{X}$, the discontinuities introduced by $\zeta_a$ at $t_{y,a}$ prevent $f_{h_a}$ from complying with individual fairness; and

(iii) if $p_a \neq p_{a'}$ then users from group $a$ cannot access the random classification odds offered to group $a'$ and vice versa.

Section 1 has thoroughly demonstrated the adverse consequences of point (i). While users are made to believe that a higher score is better, e.g., their credit rating, fixed randomisation only exhibits this behaviour at the thresholds. Refer back to Figure 1, which shows that despite there being clear evidence of *white* users with a credit score of 50 being more

likely to repay their loan than *white* candidates whose credit score is 25, both are equally likely (but not guaranteed) to receive a loan.

DEFINITION 3.1 (CLASSIFICATION ODDS DISTANCE). *Given a decision function $h_a : \mathcal{R} \mapsto \mathcal{Y}$ such that $h_a(r) \sim B(1, \zeta_a(r))$, we define the corresponding distance metric $d_{\mathcal{Y}} : I \times I \mapsto I$ such that*

$$d_{\mathcal{Y}}(h_a(r_1), h_a(r_2)) = |\zeta_a(r_1) - \zeta_a(r_2)| \qquad \forall r_1, r_2 \in \mathcal{R} .$$

*Using Equation 2, the distance is the difference in odds of positive classification between two scores.*

Point (ii) concerns the classification behaviour around the thresholds $t_{y,a}$ and fixed randomisation parameter $p_a$, which create discontinuities in odds for the final decision function $f_{h_a}$. To demonstrate this we use Definition 3.1, which specifies a distance metric on the classification odds. Lipschitz conditions scale across compositions [17], such that

$$d_{\mathcal{Y}}(h_a(g(\mathbf{x}_1)), h_a(g(\mathbf{x}_2))) \leq L_{\mathcal{R}} d_{\mathcal{R}}(g(\mathbf{x}_1), g(\mathbf{x}_2)) \leq L_{\mathcal{R}} L_{\mathcal{X}} d_{\mathcal{X}}(\mathbf{x}_1, \mathbf{x}_2) \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X} .$$

Issues arise around the thresholds. Take

$$r_1 = \lim_{z \to t_{y,a}^+} z \quad \text{and} \quad r_2 = \lim_{z \to t_{y,a}^-} z ,$$

so two scores approach a threshold from different sides. In such a case, from Equation 2, we have that

$$d_{\mathcal{Y}}(h_a(r_1), h_a(r_2)) = p_a \quad \text{or} \quad d_{\mathcal{Y}}(h_a(r_1), h_a(r_2)) = (1 - p_a) ,$$

and $d_{\mathcal{R}}(r_1, r_2) \to 0$. Therefore

$$\frac{d_{\mathcal{Y}}(h_a(r_1), h_a(r_2))}{d_{\mathcal{R}}(r_1, r_2)} \to \infty \quad \text{and} \quad \frac{d_{\mathcal{Y}}(h_a(r_1), h_a(r_2))}{d_{\mathcal{R}}(r_1, r_2)} \leq L_{\mathcal{R}} , \tag{3}$$

and thus $L_{\mathcal{R}}$ must be very large. As $r_1$ approaches $t_{y,a}$ from one side and $r_2$ from the other, $h_a$ is clearly not locally Lipschitz continuous since $d_{\mathcal{R}} \to 0$ but $d_{\mathcal{Y}} \to p_a$ or $(1-p_a)$, one of which is always above 0. In theory, $g$ could be crafted such that it cannot map individuals to values around the thresholds, however this would introduce discontinuities to $g$ and thus invalidate the Lipschitz condition. In this scenario, assuming $g$ satisfies the individual fairness constraint defined in Equation 1, $f_{h_a}$ must ultimately violate such an individual fairness constraint at the thresholds when fixed randomisation is employed. Fixed randomisation can therefore be seen as a *step function* – see Figure 1 – which is not uniformly continuous on any interval that contains $t_{y,a}$ [16]. Small changes can occur for a variety of reasons, e.g., a lack of instrumentation precision [32] or noise due to human error [34], and thus we argue that small changes should never dramatically change an individual's odds.

DEFINITION 3.2 (EQUALISED INDIVIDUAL ODDS). *Given a probabilistic classifier $f_a : \mathcal{X}_a \mapsto \mathcal{Y}$, where $\mathcal{X}_a \subseteq \mathcal{X}|A = a$, defined by the probability curve $\zeta_a : \mathcal{R} \mapsto I_a \subseteq I$ such that $h_a(r) \sim B(1, \zeta_a(r))$ and $f_a = h_a \circ g \circ \cdots$, $f_a$ satisfies individual odds iff*

$$\exists r' \in \mathcal{R} \ \ s.t. \ \ \zeta_a(r) = \zeta_{a'}(r') \quad \forall r \in \mathcal{R} \ \forall a, a' \in A \ a \neq a' .$$

*Therefore, all sub-populations in A must be capable of attaining classification odds available to all the other groups.*

Point (iii) highlights an interesting behaviour that gives rise to a novel, relatively weak, notion of fairness, which we call *individual odds* – see Definition 3.2. To satisfy this fairness criterion $\zeta_a$ does not necessarily need to be continuous but every point that it can reach must also be available to $\zeta_{a'}$, so effectively we require $I_a \equiv I_{a'}$. Violating this constraint implies that there exists a subset of users from the $A = a$ sub-population that can never be treated the same as a

portion of individuals from the $A = a'$ group and vice versa. Whenever $p_a \neq p_{a'}$, the individual odds criterion is clearly not satisfied for fixed randomisation. This definition of fairness bridges the, thus far somewhat separate, concepts of individual and group fairness as it considers the treatment of individual users in view of their assignment to distinct sub-populations determined by the protected attribute $A$.

## 4   CONSTRUCTING CURVES FOR PREFERENTIAL RANDOMISATION

Under these conditions, assuring group and individual fairness is equivalent to searching for solutions that are *continuous* and *smooth*, with a well-defined limit on $L_{\mathcal{R}}L_{\mathcal{X}}$, which also satisfy Definition 2.1. We therefore must find a combination of the group thresholds ($t_{0,a}$ and $t_{1,a}$) and a curve between them that satisfies individual as well as group fairness.

### 4.1   Defining Solution Behaviour

There are potentially infinite curves that satisfy the aforementioned conditions. In order to decrease the size of the solution space, we can impose further restrictions on the expected behaviour of the solution and parameterisation thereof. Where $h_a$ follows *fixed randomisation*, we define *preferential randomisation* as $z_a(r) \sim B\big(1, \phi_a(r)\big)$ to distinguish between the two; therefore, $f_{z_a} = z_a \circ g$ and

$$\mathbb{P}\{f_{z_a}(r) = 1 | A = a, R = r\} = \phi_a(r) \,.$$

We expect preferential randomisation to behave as follows:

**Monotonicity**  Larger values of $r$ should entail equal or higher chances of positive classification as argued by point (i) in Section 3, i.e.,

$$\phi_a'(r) \geq 0 \quad \forall r \in \mathcal{R} \,.$$

**Continuity at boundaries**  The solution should avoid sudden jumps in probability at the thresholds $t_{y,a}$ to satisfy point (ii), i.e.,

$$\phi_a(t_{y,a}) = y \,.$$

**Continuity for interval space**  The curve that maps $\mathcal{R}$ to the classification probability must be well-defined at all points in $\mathcal{R}$ in compliance with point (iii). If $\tilde{r}$ is any fixed point in $\mathcal{R}$,

$$\lim_{r \to \tilde{r}^+} \phi_a(r) = \lim_{r \to \tilde{r}^-} \phi_a(r) \quad \forall \tilde{r} \in \mathcal{R} \;,$$

where $r \to \tilde{r}^+$ is $r$ approaching $\tilde{r}$ from above and $r \to \tilde{r}^-$ is $r$ approaching $\tilde{r}$ from below.

Monotonicity between the thresholds guarantees that higher scores are treated better; continuity within the interval ensures that the Lipschitz constant does not explode at the thresholds – see Figure 3 for an example. Because no score can be outside of the $[\min(\mathcal{R}), \max(\mathcal{R})] = [\mathcal{R}_\alpha, \mathcal{R}_\omega]$ range, the output of $\phi_a$ does not need to span the entire probability range $[0, 1]$ if the thresholds are fixed at the extremes, i.e., $t_{0,a} = \mathcal{R}_\alpha \implies \phi_a(t_{0,a}) \geq 0$ or $t_{1,a} = \mathcal{R}_\omega \implies \phi_a(t_{1,a}) \leq 1$. This is especially important when we can only access the final decisions $\mathcal{Y}$ as opposed to the scores $\mathcal{R}$, i.e., $g$ is a crisp classifier $g : \mathcal{X} \mapsto \mathcal{Y}$, in which case we require the ability to randomise the crisp predictions. With these constraints we can satisfy the requirements outlined in Section 3.

### 4.2   Viable Solutions from Linear Systems

Even with these constraints, the number of curves between each combination of thresholds that constitute viable solutions is still infinite. We therefore further constrict the solution space to piece-wise polynomials parameterised

only by $t_{y,a}$ and $p_a$. We assume each solution takes the form

$$\psi_{a,y}(r) = v_y + b_y r + c_y r^2 + \cdots ,$$

and so for $\tau_a = t_{0,a} + (1 - p_a)(t_{1,a} - t_{0,a})$,

$$\phi_a(r) = \begin{cases} 0 & \text{if } r < t_{0,a} \\ \psi_{a,0}(r) & \text{if } r \in [t_{0,a}, \tau_a) \\ \psi_{a,1}(r) & \text{if } r \in [\tau_a, t_{1,a}) \\ 1 & \text{if } r \geq t_{1,a} \end{cases} . \tag{4}$$

We choose this particular point of connection ($\tau_a$) because it ensures that all solutions (including $\zeta_a$) follow

$$\int_{t_{0,a}}^{t_{1,a}} \phi_a(r) dr = \int_{t_{0,a}}^{t_{1,a}} \zeta_a(r) dr \implies \int_{\mathcal{R}} \phi_a(r) dr = \int_{\mathcal{R}} \zeta_a(r) dr .$$

This property guarantees that curves parameterised by the same thresholds and probabilities are comparable as they yield the same average probability between $t_{0,a}$ and $t_{1,a}$. The only difference between such solutions is their smoothness and continuity (see Appendix B for the proof). Finding families of closed-form solutions is achieved by using the continuity and monotonic constraints, with the addition of smoothness constraints as the order of the polynomial increases, and solving a full-rank linear system $M\mathbf{x} = \mathbf{b}$ (refer to Appendix C for details). Here, we consider four candidate curves:

**linear form**

$$\psi_{a,0}(r) = \frac{p_a(r-1)}{1 - p_a} \qquad\qquad \psi_{a,1}(r) = \frac{(1 - p_a)(r-1)}{p_a} + 1 ,$$

**quadratic form**

$$\psi_{a,0}(r) = \frac{p_a r^2}{(p_a - 1)^2} \qquad\qquad \psi_{a,1}(r) = \frac{p_a^2 + p_a - 1}{p_a^2} - \frac{2(p_a - 1)}{p_a^2} r + \frac{(p_a - 1)}{p_a^2} r^2 ,$$

**cubic form**

$$\psi_{a,0}(r) = \frac{3 p_a r^2}{p_a^2 - 2p_a + 1} + \frac{2 p_a r^3}{(p_a - 1)^3}$$

$$\psi_{a,1}(r) = \frac{p_a^3 + 3p_a^2 - 3p_a + 2}{p_a^3} + \frac{6(p_a^2 - 2p_a + 1)}{p_a^3} r + \frac{3(p_a^2 - 3p_a + 2)}{p_a^3} r^2 + \frac{2(p_a - 1)}{p_a^3} r^3 , \text{ and}$$

**4$^{\text{th}}$ order polynomial form**

$$\psi_{a,y}(r) = (30 p_a - 12) r^2 + (-60 p_a + 28) r^3 + (30 p_a - 15) r^4 . \tag{5}$$

Their derivations are given in Appendix C. Note that the 4$^{\text{th}}$ order polynomial (Equation 5) is not monotonic if $p_a \notin [\frac{2}{5}, \frac{3}{5}]$.

### 4.3 Validating Individual Fairness

If $g$ is individually fair from the outset, validating that a given solution satisfies the individual fairness constraint is straight forward. From Definition 3.1 and Equation 3,

$$d_y\big(z_a(r_1), z_a(r_2)\big) \leq L_{\mathcal{R}} d_{\mathcal{R}}(r_1, r_2) \quad \implies \quad \frac{|\phi_a(r_1) - \phi_a(r_2)|}{|r_1 - r_2|} \leq L_{\mathcal{R}} .$$
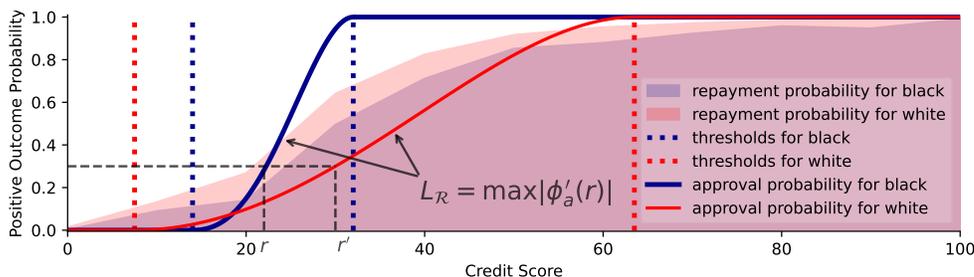
Fig. 3. Two-threshold preferential randomisation with smoothness constraints applied to probabilities (y-axis) output by a loan repayment classifier built upon credit scores (x-axis). It satisfies equalised odds for the protected attribute *race* (black and white) by using it to assign approval probabilities. This solution has *no discontinuities* – satisfying individual odds (Definition 3.2, see $r$ and $r'$ for an example) and being $L_\mathcal{R}$ Lipschitz-continuous (Equation 1) – and offers predictive performance marginally better than the fixed randomisation method shown in Figure 1.

Taking the limit $r_1 \to r_2$, we get the definition of a derivative. Therefore, we can calculate $L_\mathcal{R}$ by considering the maximum derivative $L_\mathcal{R} = \max\left(|\phi'_a(r)|\right)$ $\forall r \in [t_{0,a}, t_{1,a}]$ (proof in Appendix D). Due to the definitions of each $\phi_a$, the maximum value of $\phi'_a$ on $\mathcal{R}$ is always either at the thresholds or at the connection point, with the exception of the 4$^{\text{th}}$ order polynomial for which $L_\mathcal{R}$ is where $\phi''_a(r) = 0$ for $r \in (t_{0,a}, t_{1,a})$, thus $L_\mathcal{R}$ is always known. Finding an optimal solution is therefore a case of identifying values of $t_{y,a}$ and $p_a$ for $\phi_a$ that satisfy Definition 2.1 such that $L_\mathcal{R} L_\mathcal{X}$ is well-defined. While $L_\mathcal{R}$ is not guaranteed to be small, it is guaranteed to be finite. Taking the limit $p_a \to 1$ or 0, $t_{0,a} \to t_{1,a}$, or $t_{1,a} \to t_{0,a}$, then $L_\mathcal{R} \to \infty$, which is synonymous with using a single threshold, hence invalidating equalised odds.

## 5 CASE STUDIES

Here we apply the method of preferential randomisation to two case studies: credit scoring for loan allocation (CreditRisk) and risk of recidivism (COMPAS). Source code for all the studies is available online[1].

### 5.1 CreditRisk Case Study

To facilitate a direct comparison, we apply our method to the case study conducted by Hardt et al. [22]. Credit scores are often used to determine whether an individual should receive a loan or mortgage, to calculate interest rates and credit limits, and even to conduct background check on tenants [21, 29]. The scoring function $g$ – which calculates credit scores on input space $\mathcal{X}$ – operates as a black box (see Section 3), therefore we only observe the scores $R$ and cannot access $\mathcal{X}$ or $g$.

The input space may contain attributes influenced by cultural background (i.e., related to race), possibly causing the joint distribution of $R$ and $Y$ to differ between sub-populations $A$. The CreditRisk data set captures the credit score's ability to predict defaulting on a loan (i.e., failing to repay it) for 90 days or more. The data show that as credit score increases, the likelihood of defaulting decreases (shown in Figure 8 given in Appendix E) . The rate of these changes, however, is correlated with *race*. Therefore, when a single threshold for each sub-population is optimised for maximum accuracy, the equalised odds (Definition 2.1) becomes 0.28; we should strive for this fairness metric to be as close to 0 as possible.

---

[1]https://github.com/teddyzander/McGIF

| | white | | | black | | | Asian | | |
|---|---|---|---|---|---|---|---|---|---|
| | acc (%) | EO | $L_{\mathcal{R}}$ | acc (%) | EO | $L_{\mathcal{R}}$ | acc (%) | EO | $L_{\mathcal{R}}$ |
| fixed | 82.424 | 0.875 | $\infty$ | 81.483 | 0.382 | $\infty$ | 82.540 | **0.175** | $\infty$ |
| linear | **82.435** | 0.337 | 0.044 | 81.480 | 0.745 | 0.222 | 82.544 | 0.577 | 0.148 |
| quadratic | 82.440 | 1.336 | 0.046 | 81.486 | 0.558 | 0.416 | 82.543 | 0.347 | 0.115 |
| cubic | 82.430 | **0.241** | 0.091 | **81.487** | **0.103** | 0.338 | 82.542 | 0.324 | 0.265 |
| $4^{th}$ order | 82.432 | 0.428 | **0.027** | 81.482 | 0.569 | **0.092** | **82.545** | 0.554 | **0.064** |

Table 1. Accuracy (acc) as a percentage, equalised odds (EO) to the order of $\times 10^{-4}$, and Lipschitz constant ($L_{\mathcal{R}}$) per method for each value of the protected attribute *race* in the CreditRisk loan repayment prediction task. The Hispanic group is not shown as it uses a single threshold ($t_{y,a} = 30$) due to having the lowest ROC curve at the optimum, thus acting as the baseline for other *races*.

We overcome this by using different thresholds and probabilities (specified in Table 3 given in Appendix F) achieved with a set of curves with differing smoothness constraints. These curves honour the "higher credit score leads to higher repayment probability" dependency encoded in the underlying data. Referring back to the example introduced in Section 1, we can see from Figures 3 and 4 that the *white* user with a credit score of 49.5 is now 2.6–5.25 times more likely to receive a loan than the *white* user with a credit score of 25, depending on which continuous solution is chosen. The results – reported in Table 1 – show that the difference in accuracy and equalised odds between fixed randomisation and preferential randomisation is negligible (a change of +0.016 and −0.000634 respectively). The method additionally improves individual fairness by the Lipschitz constant on $\phi_a$ and through satisfying Definition 3.2 (individual odds). Preferential randomisation can therefore be used to guarantee group and individual fairness through the notions of *equalised odds* and *individual odds*, and this encourages users to engage with the scoring model.
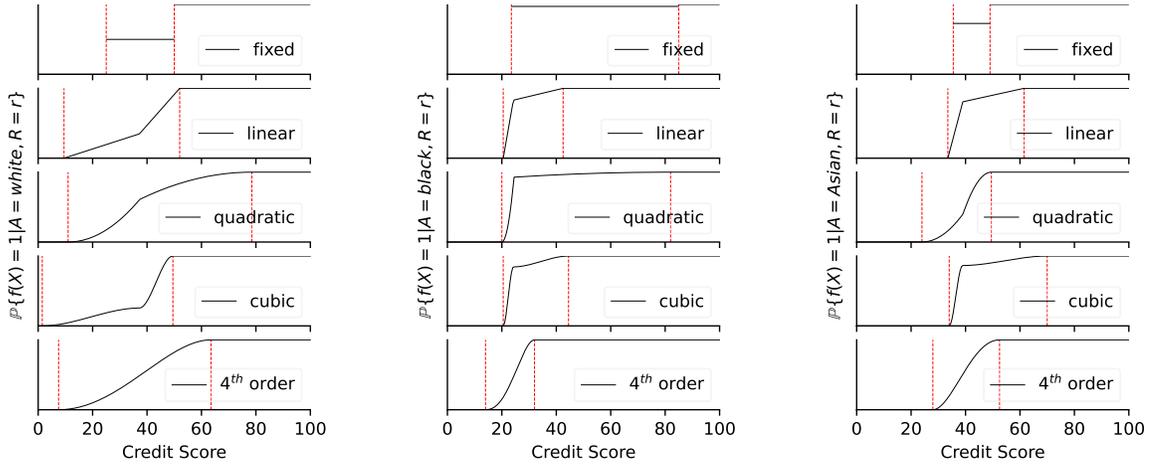


Fig. 4. Probability curves corresponding to the results reported in Table 1. All solutions have comparable accuracy and satisfy equalised odds but yield a different Lipschitz constant $L_{\mathcal{R}}$. The Hispanic group is omitted as it uses a single threshold $t_{y,a} = 30$ (refer to Table 3 given in Appendix F).
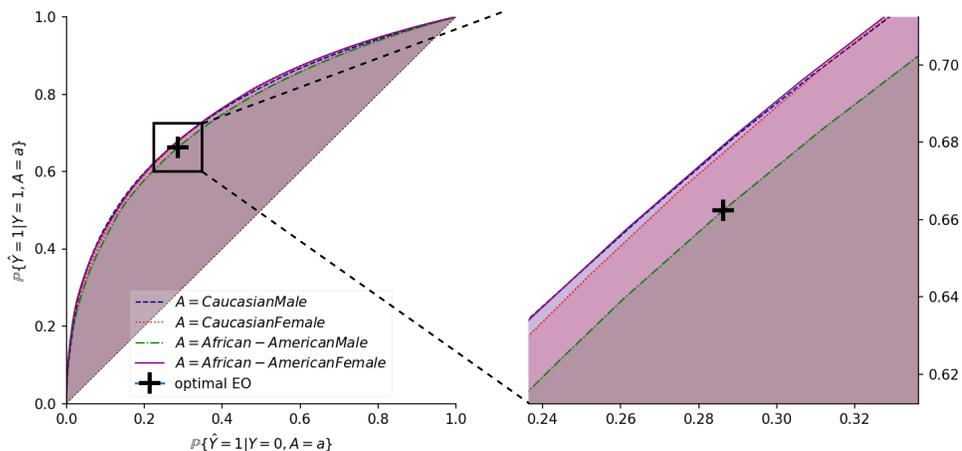
Fig. 5. ROC curves for the COMPAS data set. The coloured regions indicate areas accessible to each group. (Refer to Figure 2 for more details.)

## 5.2 COMPAS Case Study

The COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) software[2] is a commercial tool used across multiple U.S. states to analyse and predict a defendant's behaviour if released on bail. The software output can be considered by judges during sentencing, albeit such a practice must be disclosed. Specifically, COMPAS offers three insights: (1) likelihood of general recidivism (re-offending); (2) likelihood of violent recidivism (committing a violent crime); and (3) likelihood of failing to appear in court (pretrial flight risk). Here, we focus on the risk of re-offending using the raw COMPAS scores available in the ProPublica data set[3] [2]. The COMPAS algorithm uses characteristics such as criminal history, known associates, drug involvement and indicators of juvenile delinquency in order to calculate a score $r$, where a higher score corresponds to a higher likelihood of recidivism. As is the case with CreditRisk (Section 5.1), the scoring algorithm used by the COMPAS software is proprietary. Given its high stakes nature, it is important to understand the predictive behaviour of this tool since its social situatedness captured by the (protected) data features – which are translated into the score $r$ – may yield biased results [35] as shown in Figure 5.

To this end, we define $A$ as the Cartesian product of two sensitive attributes – *sex* $A_1$ = {male, female} and *race* $A_2$ = {Caucasian, African-American} – found in the COMPAS data set, such that $A = A_1 \times A_2$ and so $|A| = 4$. Additionally, normalised COMPAS scores for a population of interest are denoted with $R$; the ground truth label for each score in $R$ is given by $Y$, where 1 corresponds to individuals who committed an offence in a two-year time window; and $\widehat{Y}$ captures crisp predictions, with 1 indicating high risk (of recidivism). Studying the link between the scores and labels provided by the COMPAS data set – refer to Figure 9 given in Appendix E – indicates that for most values of $r$ across all groups encoded by $A$ if $r_1 > r_2$, then $\mathbb{P}\{Y = 1|R = r_1\} \geq \mathbb{P}\{Y = 1|R = r_2\}$. Therefore, we are in a good position to use the monotonic probability functions proposed in this paper to build the final classifier.

| | Caucasian male | | | Caucasian female | | | African-American female | | |
|---|---|---|---|---|---|---|---|---|---|
| | acc (%) | EO | $L_\mathcal{R}$ | acc (%) | EO | $L_\mathcal{R}$ | acc (%) | EO | $L_\mathcal{R}$ |
| fixed | **67.400** | 2.412 | $\infty$ | **67.595** | 1.413 | $\infty$ | 67.549 | 1.625 | $\infty$ |
| linear | 67.384 | **0.632** | 0.181 | 67.574 | 1.190 | **0.073** | 67.555 | **0.488** | 0.109 |
| quadratic | 67.386 | 0.640 | 0.254 | 67.594 | 1.723 | 0.129 | 67.548 | 1.259 | 0.214 |
| cubic | 67.395 | 0.709 | 0.254 | 67.589 | **0.203** | 0.513 | 67.555 | 2.349 | 0.427 |
| $4^{\text{th}}$ order | 67.380 | 1.145 | **0.075** | 67.576 | 2.431 | 0.080 | **67.574** | 1.863 | **0.072** |

Table 2. Accuracy (acc) as a percentage, equalised odds (EO) to the order of $\times 10^{-4}$, and Lipschitz constant ($L_\mathcal{R}$) per method for each value of the Cartesian product of the protected attributes *race* and *sex* in the COMPAS prediction task. The African-American male group is not shown as it uses a single threshold ($t_{y,a} = 48$) due to having the lowest ROC curve at the optimum, thus acting as the baseline for other groups.

Given the aforementioned relationship, it is in the public's (and judicial system's) best interest to always increase the probability of classifying an individual as high-risk when the score $r$ increases. However, fixed randomisation does not allow for this. For example, under fixed randomisation a *Caucasian male* with a COMPAS score in the $[24, 41)$ range has an 11.6% chance of being classified as high-risk (see Table 4 in Appendix F); nonetheless, a *Caucasian male* at the top of this score range is almost twice as likely to commit an offence as a *Caucasian male* with a score at the low end of this range. Therefore, fixed randomisation is unfair on three fronts: (1) *Caucasian males* with scores in the lower range of $[24, 42)$ are treated the same as *Caucasian males* with scores in the higher range of this interval; (2) higher risk individuals are not labeled as such despite their scores indicating so; and (3) individuals whose outcome is randomised are never offered the same odds as members of other protected groups (in violation of Definition 3.2). Notably, these arguments apply to all groups in the protected attribute $A$ and not only *Caucasian males*. Small changes in score having a large impact on odds can have very real effects on individuals – see the case of Mr. Rodriguez, whose
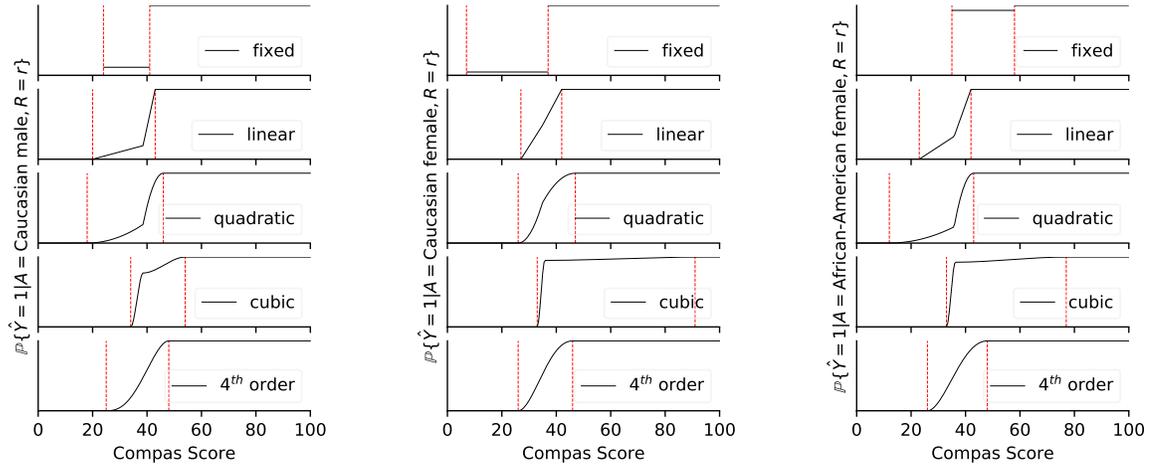


Fig. 6. Probability curves corresponding to the results reported in Table 2. All solutions have comparable accuracy and satisfy equalised odds but yield a different Lipschitz constant $L_\mathcal{R}$. The African-American male group is omitted as it uses a single threshold $t_{y,a} = 48$ (refer to Table 4 given in Appendix F).

analysis contained an error that caused his parole to be incorrectly denied [42] – thus we argue that small changes should not dramatically change the chances of classification.

The between-group equalised odds measure when we maximise accuracy separately for each group is 0.148. Mirroring Section 5.1, we apply our method to the COMPAS model in order to reduce the equalised odds disparity without breaking individual fairness. We therefore seek to calibrate the model with a combination of thresholds and probabilities that parameterise the continuous curves (defined in Section 4.2) using nothing but the joint distributions of $(R, A, Y)$. We then compare the continuous solutions to the step function solution (fixed randomisation) defined in Equation 2. The results reported in Table 2 and Figure 6 show that continuous curves can be used to simultaneously satisfy equalised odds, individual fairness and individual odds. Since low-scoring individuals are less likely to be classified as high-risk, defendants have an incentive to engage in behaviour that actively lowers their COMPAS score. Furthermore, public safety is prioritised more effectively since individuals with a measurably higher probability of recidivism are given higher odds of being classified as high-risk.

## 6 CONCLUSION AND FUTURE WORK

In this work we demonstrated how using fixed randomisation to guarantee group fairness may be detrimental to both the owners and users of a predictive model. Users with higher scores should be more likely to receive a better outcome – a property that may be lost when enforcing group fairness. Ensuring this behaviour also allows the owners to preserve predictive performance and transparency of the automated decision-making process. By using the method proposed in this paper – which relies on monotonic and continuous curves – we can guarantee these properties. Our approach rewards building accurate scoring functions and adheres to the notion of individual fairness from the perspective of function composition. Importantly, the burden of accurate classification remains the sole responsibility of the model owner since our method forces all individuals to rely on the equalised odds measure of the worst-performing sub-population. This allocation of responsibility is desirable as owners can choose to invest in better predictors, data or scoring functions, whereas users in under-performing groups lack this agency.

Notably, our case study shows that there can exist multiple solutions that simultaneously satisfy equalised odds and individual fairness, which can be linked to *model multiplicity* [36]. When equalised odds, individual fairness and accuracy are comparable between groups, we can choose to discriminate the solutions based on other criteria. Future work will explore this aspect of our curves; specifically, we will consider:

(1) the most robust curve for each group [26];
(2) curves such that $L_{\mathcal{R}}$ is closest between groups;
(3) the smoothest curves;
(4) curves that subject the fewest individuals to random outcomes, for example, $\min|t_{1,a} - t_{0,a}| \ \forall a \in A$ ; and
(5) curves that subject equal number of individuals to random outcomes between groups, e.g., $\min \sum_{\forall a \in A} \left( |t_{1,a} - t_{0,a}| - |t_{1,a'} - t_{0,a'}| \right)$ where $a \neq a'$.

# REFERENCES

[1] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. 2018. A reductions approach to fair classification. In *International conference on machine learning*. 60–69.

[2] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias. *ProPublica* 23 (May 2016).

[3] Pranjal Awasthi, Matthäus Kleindessner, and Jamie Morgenstern. 2020. Equalized odds postprocessing under imperfect group information. In *Proceedings of the 23$^{rd}$ International Conference on Artificial Intelligence and Statistics*, Vol. 108. PMLR, 1770–1780.

[4] Pablo Barceló, Mikaël Monet, Jorge Pérez, and Bernardo Subercaseaux. 2020. Model Interpretability through the lens of computational complexity. In *Advances in Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc., 15487–15498.

[5] Joachim Baumann, Anikó Hannák, and Christoph Heitz. 2022. Enforcing Group Fairness in Algorithmic Decision Making: Utility Maximization Under Sufficiency. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 2315–2326.

[6] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. 2006. Manifold Regularization: A Geometric Framework for Learning from Labeled and Unlabeled Examples. *Journal of Machine Learning Research* 7, 85 (2006), 2399–2434.

[7] Reuben Binns. 2020. On the Apparent Conflict between Individual and Group Fairness. In *Proceedings of the 2020 ACM Conference on Fairness, Accountability, and Transparency*. 514–524.

[8] Sumon Biswas and Hridesh Rajan. 2021. Fair Preprocessing: Towards Understanding Compositional Fairness of Data Transformers in Machine Learning Pipeline. In *Proceedings of the 29$^{th}$ ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 981–993.

[9] Alessandro Castelnovo, Riccardo Crupi, Greta Greco, Daniele Regoli, Ilaria Penco, and Andrea Cosentini. 2022. A clarification of the nuances in the fairness metrics landscape. *Scientific Reports* 12 (2022).

[10] Irene Chen, Fredrik D Johansson, and David Sontag. 2018. Why Is My Classifier Discriminatory?. In *Advances in Neural Information Processing Systems*, Vol. 31. Curran Associates, Inc.

[11] Yan Chen, Christopher Mahoney, Isabella Grasso, Esma Wali, Abigail Matthews, Thomas Middleton, Mariama Njie, and Jeanna Matthews. 2021. Gender Bias and Under-representation in Natural Language Processing Across Human Languages. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 24–34.

[12] Mihai Cimpoeşu, Andrei Sucilă, and Henri Luchian. 2013. A Statistical Binary Classifier: Probabilistic Vector Machine. In *Progress in Artificial Intelligence: 16$^{th}$ Portuguese Conference on Artificial Intelligence, EPIA 2013*. Springer, 211–222.

[13] Krzysztof Dembczyński, Weiwei Cheng, and Eyke Hüllermeier. 2010. Bayes Optimal Multilabel Classification via Probabilistic Classifier Chains. In *Proceedings of the 27$^{th}$ International Conference on International Conference on Machine Learning (ICML'10)*. 279–286.

[14] Sanghamitra Dutta, Dennis Wei, Hazar Yueksel, Pin-Yu Chen, Sijia Liu, and Kush Varshney. 2020. Is There a Trade-off Between Fairness and Accuracy? A Perspective Using Mismatched Hypothesis Testing. In *Proceedings of the 37$^{th}$ International Conference on Machine Learning*, Vol. 119. PMLR, 2803–2813.

[15] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through Awareness. In *Proceedings of the 3$^{rd}$ Innovations in Theoretical Computer Science Conference*. 214–226.

[16] Donald Estep. 2002. *Practical Analysis in One Variable*. Springer. 83–87 pages.

[17] Herbert Federer. 1996. *Geometric Measure Theory*. Springer Berlin Heidelberg.

[18] Tilmann Gneiting. 2011. Making and Evaluating Point Forecasts. *J. Amer. Statist. Assoc.* 106, 494 (2011), 746–762.

[19] Tilmann Gneiting and Adrian E Raftery. 2007. Strictly Proper Scoring Rules, Prediction, and Estimation. *J. Amer. Statist. Assoc.* 102, 477 (2007), 359–378.

[20] Nina Grgic-Hlaca, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. 2016. The case for process fairness in learning: Feature selection for fair decision making. In *NIPS symposium on machine learning and the law*.

[21] Andrew Hanson, Zackary Hawley, Hal Martin, and Bo Liu. 2016. Discrimination in mortgage lending: Evidence from a correspondence experiment. *Journal of Urban Economics* 92 (2016), 48–65.

[22] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of Opportunity in Supervised Learning. In *Advances in Neural Information Processing Systems*, Vol. 29. Curran Associates, Inc.

[23] Tatiana Homonoff, Rourke O'Brien, and Abigail B Sussman. 2021. Does Knowing Your FICO Score Change Financial Behavior? Evidence from a Field Experiment with Student Loan Borrowers. *The Review of Economics and Statistics* 103, 2 (2021), 236–250.

[24] Krzysztof Kil, Radosław Ciukaj, and Justyna Chrzanowska. 2021. Scoring Models and Credit Risk: The Case of Cooperative Banks in Poland. *Risks* 9, 7 (2021).

[25] Pranay K Lohia, Karthikeyan Natesan Ramamurthy, Manish Bhide, Diptikalyan Saha, Kush R Varshney, and Ruchir Puri. 2019. Bias Mitigation Post-processing for Individual and Group Fairness. In *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2847–2851.

[26] Xinsong Ma, Zekai Wang, and Weiwei Liu. 2022. On the Tradeoff Between Robustness and Fairness. In *Advances in Neural Information Processing Systems*.

[27] Anton Markov, Zinaida Seleznyova, and Victor Lapshin. 2022. Credit scoring methods: Latest trends and points to consider. *The Journal of Finance and Data Science* 8 (2022), 180–201.

[28] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A Survey on Bias and Fairness in Machine Learning. *Comput. Surveys* 54, 6 (2021).

[29] Anthony Pennington-Cross. 2003. Credit history and the performance of prime and nonprime mortgages. *The Journal of Real Estate Finance and Economics* 27, 3 (2003), 279–301.

[30] Felix Petersen, Debarghya Mukherjee, Yuekai Sun, and Mikhail Yurochkin. 2021. Post-processing for Individual Fairness. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 25944–25955.

[31] Alvin Rajkomar, Michaela Hardt, Michael D Howell, Greg Corrado, and Marshall H Chin. 2018. Ensuring Fairness in Machine Learning to Advance Health Equity. *Annals of Internal Medicine* 169, 12 (2018), 866–872.

[32] Manon Romain, Adrien Senecat, Soizic Pénicaud, Gabriel Geiger, and Justin-Casimir Braun. 2023. How We Investigated France's Mass Profiling Machine. https://www.lighthousereports.com/methodology/how-we-investigated-frances-mass-profiling-machine/

[33] Yaniv Romano, Stephen Bates, and Emmanuel Candes. 2020. Achieving Equalized Odds by Resampling Sensitive Attributes. In *Advances in Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc., 361–371.

[34] Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence* 1, 5 (2019), 206–215.

[35] Cynthia Rudin, Caroline Wang, and Beau Coker. 2020. The age of secrecy and unfairness in recidivism prediction. *Harvard Data Science Review* 2, 1 (2020), 1.

[36] Kacper Sokol, Meelis Kull, Jeffrey Chan, and Flora Dilys Salim. 2022. Cross-model Fairness: Empirical Study of Fairness and Ethics Under Model Multiplicity. *arXiv preprint arXiv:2203.07139* (2022).

[37] Laura Spring, Diana Robillard, Lorrie Gehlbach, and Tiffany A Moore Simas. 2011. Impact of pass/fail grading on medical students' well-being and academic outcomes. *Medical Education* 45, 9 (2011), 867–877.

[38] U.S. Federal Reserve. 2007. Report to the congress on credit scoring and its effects on the availability and affordability of credit. *Board of Governors of the Federal Reserve System* (2007).

[39] Robin Vogel, Aurélien Bellet, and Stephan Clémençon. 2021. Learning Fair Scoring Functions: Bipartite Ranking under ROC-based Fairness Constraints. In *International conference on artificial intelligence and statistics*. PMLR, 784–792.

[40] Mingyang Wan, Daochen Zha, Ninghao Liu, and Na Zou. 2022. In-processing Modeling Techniques for Machine Learning Fairness: A Survey. *ACM Transactions on Knowledge Discovery from Data* (2022).

[41] Yuyan Wang, Xuezhi Wang, Alex Beutel, Flavien Prost, Jilin Chen, and Ed H Chi. 2021. Understanding and Improving Fairness–Accuracy Trade-offs in Multi-task Learning. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining* (2021).

[42] Rebecca Wexler. 2017. When a computer program keeps you in jail: How computers are harming criminal justice. *New York Times* 13 (2017).

[43] Mark Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gaby Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Olavo Bonino da Silva Santos, Philip Bourne, Jildau Bouwman, Anthony Brookes, Tim Clark, Merce Crosas, Ingrid Dillo, Olivier Dumon, Scott Edmunds, Chris Evelo, Richard Finkers, and Barend Mons. 2016. The FAIR guiding principles for scientific data management and stewardship. *Scientific Data* 3 (2016).

[44] Sang Eun Woo, James M LeBreton, Melissa G Keith, and Louis Tay. 2023. Bias, Fairness, and Validity in Graduate-school Admissions: A Psychometric Perspective. *Perspectives on Psychological Science* 18, 1 (2023), 3–31.

[45] Han Xu, Xiaorui Liu, Yaxin Li, Anil Jain, and Jiliang Tang. 2021. To be robust or to be fair: Towards fairness in adversarial training. In *Proceedings of the 38th International Conference on Machine Learning*. 11492–11501.

[46] Chaojian Yu, Bo Han, Li Shen, Jun Yu, Chen Gong, Mingming Gong, and Tongliang Liu. 2022. Understanding Robust Overfitting of Adversarial Training and Beyond. In *Proceedings of the 39th International Conference on Machine Learning*, Vol. 162. PMLR, 25595–25610.

## A    EXAMPLE DISTANCE FUNCTIONS

If $\mathcal{M}$ is a metric space and $a, b, c \in \mathcal{M}$, then $d_{\mathcal{M}} : \mathcal{M} \times \mathcal{M} \mapsto \mathbb{R}^+$ and the following hold:

- $d_{\mathcal{M}}(a, a) = 0$ – the distance between a point and itself is 0;
- if $a \neq b$, $d_{\mathcal{M}}(a, b) > 0$ – the distance between two different points is strictly greater than 0;
- $d_{\mathcal{M}}(a, b) = d_{\mathcal{M}}(b, a)$ – the distance between two different points $a$ and $b$ is equal to the distance between $b$ and $a$; and
- $d_{\mathcal{M}}(a, c) \leq d_{\mathcal{M}}(a, b) + d_{\mathcal{M}}(b, c)$ – the distance between any two points is equal to or less than the distance given by visiting another point on a journey between the original two points (triangle inequality).

### A.1    Euclidean Distance (Continuous Features)

The $L^2$-norm is defined as

$$||\mathbf{x}||_2 = \sqrt{\sum_{k=1}^{N} |x_k|^2} \, ,$$

and is the foundation of Euclidean distance $d_E : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ defined as

$$d_E(\mathbf{x}_1, \mathbf{x}_2) = ||\mathbf{x}_1 - \mathbf{x}_2||_2 \, ,$$

and so

$$d_E(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{\sum_{k=1}^{N} |x_{1,k} - x_{2,k}|^2} \, .$$

### A.2    Hamming Distance (Discrete Features)

The $L^1$-norm is defined as

$$||\mathbf{x}||_1 = \sum_{k=1}^{N} |x_k| \, , \tag{6}$$

and is the foundation of Hamming distance $d_H : \mathcal{X} \times \mathcal{X} \mapsto \{0, 1, \ldots, N-1, N\}$, which counts the number of features that differ between two inputs $\mathbf{x}_1$ and $\mathbf{x}_2$, and is defined as

$$d_H(\mathbf{x}_1, \mathbf{x}_2) = ||\mathbf{x}_1 \oplus \mathbf{x}_2||_1 \, ,$$

where $\oplus$ is the XOR operation. Therefore, $\mathbf{x}_1 \oplus \mathbf{x}_2$ is simply a vector of 0's and 1's such that

$$(\mathbf{x}_1 \oplus \mathbf{x}_2)_k = \begin{cases} 0 & \text{if } x_{1,k} = x_{2,k} \\ 1 & \text{if } x_{1,k} \neq x_{2,k} \end{cases} .$$

For example,

$$\mathbf{x}_1 = \begin{bmatrix} 3 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \\ 4 \\ 2 \end{bmatrix} \quad \mathbf{x}_2 = \begin{bmatrix} 3 \\ 0 \\ 2 \\ 0 \\ 0 \\ 0 \\ 4 \\ 1 \end{bmatrix} \implies \mathbf{x}_1 \oplus \mathbf{x}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} \implies ||\mathbf{x}_1 \oplus \mathbf{x}_2||_1 = 5 \, .$$

### A.3 Gower's Distance (Mixed Continuous and Discrete Features)

Take $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ that contain both continuous (numerical) variables and discrete (categorical) variables. We then consider each variable for $k = 1, \ldots, n$. If the pair $x_{1,k}, x_{2,k}$ is continuous,

$$s_k = 1 - \frac{|x_{1,k} - x_{2,k}|}{V_n} \, ,$$

where $V_n$ is the range of the $k^{\text{th}}$ feature. Fundamentally, the second term is the normalised $L^1$-norm (defined in Equation 6) on the differences between two vectors. However, if the pair $x_{1,k}, x_{2,k}$ is discrete, we use the Iverson operation defined by

$$s_k = [x_{1,k} = x_{2,k}] = \begin{cases} 0 & \text{if } x_{1,k} \neq x_{2,k} \\ 1 & \text{if } x_{1,k} = x_{2,k} \end{cases} \, .$$

As such, a value of $s_k = 1$ for both continuous and discrete features implies that $x_{1,k} = x_{2,k}$, and $s_k = 0$ implies that $x_{1,k}$ and $x_{2,k}$ are maximally different. We put this together to get Gower's Similarity Coefficient

$$S_G(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{n} \sum_{k=1}^{n} s_k \, , \tag{7}$$

which is bounded within $[0, 1]$. However, this coefficient does not follow the axioms laid out at the beginning of this section as $S_G(a, a) = 1$. Therefore, using Equation 7 we define Gower's distance as

$$d_G(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{1 - S_G(\mathbf{x}_1, \mathbf{x}_2)} \, ,$$

which offers the behaviour expected of a distance metric.

## B   GEOMETRIC MOTIVATION FOR THE AVERAGE PROBABILITY (LINEAR CASE)

We can parameterise each set of potential solutions for each value of $A$ (i.e., protected sub-population) by only three parameters – $p_a$, $t_{0,a}$ and $t_{1,a}$ – by constraining the area under each curve as equal to $p_a(t_{1,a} - t_{0,a})$. This forces all potential solutions with the same set of parameters to have the same average probability between the thresholds.

### B.1   $\tau_a$ Proof (Point of Intersection)

Here, we discuss the details of bounding the piece-wise linear solution such that the two lines join at $\tau_a \in [t_{0,a}, t_{1,a}]$. We present a proof that this value can be easily found, and is defined only through $q_a = 1 - p_a$, $t_{0,a}$ and $t_{1,a}$.
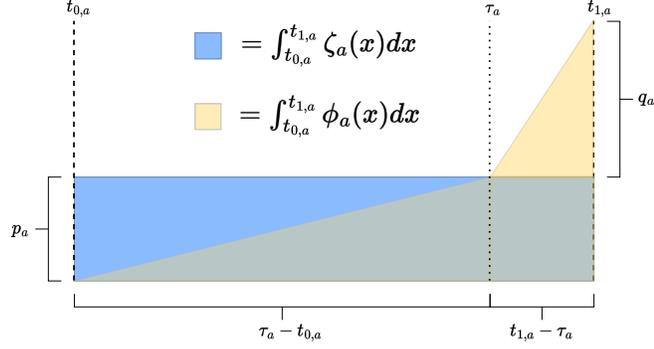
Fig. 7. Geometric interpretation of a piece-wise linear solution using thresholds and probabilities. We can see that a step function yields a rectangular area – blue space that denotes the average probability – and is defined by the probability $p_a$ and the threshold interval size $t_{1,a} - t_{0,a}$. We expect this value to be equal to the area bounded by the piece-wise linear solution $\phi(x)$ and the x-axis (yellow space), which can be decomposed into simple geometric shapes and summed up.

We begin by assuming that the solution is linear in nature and, as outlined in the paper (Section 4.1), it preserves the average probability of the step function within the interval $[t_{0,a}, t_{1,a}]$. As such, we can generate a geometric interpretation of the solution as shown in Figure 7. From here, we can see that finding $\tau_a$ is straight forward. By combining the area of a triangle equation, area of a rectangle equation, and forming an equality between equations, we get

$$p_a(t_{1,a} - t_{0,a}) = \frac{p_a}{2}(\tau_a - t_{0,a}) + \frac{q_a}{2}(t_{1,a} - \tau_a) + (t_{1,a} - \tau_a)p_a \,,$$

which can be interpreted as *Rectangle = Triangle 1 + Triangle 2 + Small Rectangle*. We can then rearrange the notation as follows:

$$\frac{p_a}{2}(2T_{1,a} - 2T_{0,a} - \tau_a + t_{0,a} - 2T_{1,a} + 2\tau_a) = \frac{q_a}{2}(t_{1,a} - \tau_a)$$

$$\frac{p_a}{2}(\tau_a - t_{0,a}) = \frac{q_a}{2}(t_{1,a} - \tau_a)$$

$$p_a(\tau_a - t_{0,a}) = q_a(t_{1,a} - \tau_a) \,.$$

Since we know that $p_a = 1 - q_a$,

$$(1 - q_a)(\tau_a - t_{0,a}) = q_a(t_{1,a} - \tau_a)$$

$$\tau_a - t_{0,a} - q_a(\tau_a - t_{0,a}) = q_a(t_{1,a} - \tau_a)$$

$$\tau_a - t_{0,a} = q_a(\tau_a - \tau_a + t_{1,a} - t_{0,a})$$

$$\tau_a = t_{0,a} + q_a(t_{1,a} - t_{0,a}) \,.$$

We therefore define $\Delta T_a = t_{1,a} - t_{0,a}$ to get the final result from Section 4.2:

$$\tau_a = t_{0,a} + q_a \Delta T_a \,. \tag{8}$$

□

## B.2 $\psi_1(x)$ and $\psi_0(x)$ Proof (Piece-wise Linear Solution)

We define the linear form of the interpolant as

$$\phi_a(x) = \begin{cases} 0 & x < t_{0,a} \\ \psi_1(x) & t_{0,a} \le x < \tau_a \\ \psi_0(x) & \tau_a \le x < t_{1,a} \\ 1 & x \ge t_{1,a} \end{cases}, \tag{9}$$

where each $\psi_n$ is linear, so

$$\psi_1(x) = vx + b \qquad \psi_0(x) = cx + d .$$

In order to derive the final form, we must assume the following conditions (continuity):

(1) $\psi_1(t_{0,a}) = 0$,

(2) $\psi_0(t_{1,a}) = 1$, and

(3) $\psi_1(\tau_a) = \psi_0(\tau_a) = p_a$.

From conditions 1 and 3 we get

$$v T_{0,a} + b = 0 \qquad v\tau_a + b = p_a .$$

From the difference of these equations we get

$$a(\tau_a - t_{0,a}) = p_a .$$

From the proof in Appendix B.1, we know that $\tau_a = t_{0,a} + q_a \Delta T_a$ (Equation 8), giving

$$v(t_{0,a} + q_a \Delta T_a - t_{0,a}) = p_a$$

$$v q_a \Delta T_a = p_a$$

$$v = \frac{p_a}{\Delta T_a q_a} .$$

Similarly, from conditions 2 and 3 we get

$$c t_{1,a} + d = 1 \qquad c\tau_a + d = p_a .$$

The difference yields

$$c(t_{1,a} - \tau_a) = 1 - p_a .$$

From the definition of $\tau_a$ and $p_a$ we get

$$c(t_{1,a} - t_{0,a} - \Delta T_a q_a) = q_a$$

$$c\Delta T_a(1 - q_a) = q_a$$

$$c\Delta T_a p_a = q_a$$

$$c = \frac{q_a}{p_a \Delta T_a} .$$

Substituting these values back into the equations yields the final parameters

$$-b = \frac{p_a}{\Delta T_a q_a} t_{0,a} \qquad\qquad -d = \frac{q_a}{p_a \Delta T_a} t_{1,a} - 1$$

$$b = -\frac{p_a T_{0,a}}{\Delta T_a q_a} \qquad\qquad d = 1 - \frac{q_a T_{1,a}}{p_a \Delta T_a} \;.$$

Putting it all back into the piece-wise equation and factorising then gives

$$\phi_a(x) = \begin{cases} 0 & x < t_{0,a} \\ \frac{p_a}{\Delta T_a q_a}(x - t_{0,a}) & t_{0,a} \le x < \tau_a \\ \frac{q_a}{\Delta T_a p_a}(x - t_{1,a}) + 1 & \tau_a \le x < t_{1,a} \\ 1 & x \ge t_{1,a} \end{cases} \qquad (10)$$

□

## B.3 Preservation of Average Probability Proof

We previously stated that solutions with the same set of parameters should always have the same average probability. For example, the linear solution $\phi_a$ preserves group fairness introduced by the step function by maintaining the same predictive behaviour (on average) in the $[t_{0,a}, t_{1,a}]$ interval. This follows directly from how we defined (1) the linear solution and (2) the point of intersection. Nonetheless, we can also prove this property directly. If $\zeta_a$ is the step function for $A = a$, we follow by stating that we require

$$\int_{-\infty}^{\infty} \zeta_a(x)dx = \int_{-\infty}^{\infty} \phi_a(x)dx \;.$$

The first thing to observe is that

$$\zeta_a(x) = \phi_a(x) \qquad \text{if } x < t_{0,a}$$

$$\zeta_a(x) = \phi_a(x) \qquad \text{if } x \ge t_{1,a}$$

from Equations 2 and 4, and so

$$\int_{-\infty}^{t_{0,a}} \zeta_a(x)dx = \int_{-\infty}^{t_{0,a}} \phi_a(x)dx$$

$$\int_{t_{1,a}}^{\infty} \zeta_a(x)dx = \int_{t_{1,a}}^{\infty} \phi_a(x)dx \;.$$

We know that the integral in the interval for the step function is

$$\int_{t_{0,a}}^{t_{1,a}} \zeta_a(x)dx = \int_{t_{0,a}}^{t_{1,a}} p_a dx$$

$$= p_a(t_{1,a} - t_{0,a})$$

$$= p_a \Delta T_a \;.$$

From Equation 9 in Appendix B.2 we know that

$$\psi_1(x) = \frac{p_a}{\Delta T_a q_a}(x - t_{0,a}) \qquad \psi_0(x) = \frac{q_a}{\Delta T_a p_a}(x - t_{1,a}) + 1 \;. \qquad (11)$$

We can decompose the integral of the piece-wise linear solution into two integrals over the interval, so using Equation 11 we get

$$\int_{t_{0,a}}^{t_{1,a}} \phi_a(x)dx = \int_{t_{0,a}}^{\tau_a} \psi_1(x)dx + \int_{\tau_a}^{t_{0,a}} \psi_0(x)dx \ . \tag{12}$$

Therefore, from Equation 10 we get

$$\int_{t_{0,a}}^{\tau_a} \psi_1(x)dx = \int_{t_{0,a}}^{\tau_a} \frac{p_a}{\Delta T_a q_a}(x - t_{0,a})dx$$
$$= \frac{p_a(\tau_a - t_{0a})^2}{2\Delta T_a q_a} \ .$$

From definition of $\Delta T_a$ and $\tau_a$ in Appendix B.1 we then get

$$\frac{p_a(\tau_a - t_{0a})^2}{2\Delta T_a q_a} = \frac{(t_{1,a} - t_{0,a})p_a q_a}{2} \ . \tag{13}$$

Also from Equation 10, we have

$$\int_{\tau_a}^{t_{1,a}} \psi_0(x)dx = \int_{\tau_a}^{t_{0,a}} \frac{q_a}{\Delta T_a p_a}(x - t_{1,a}) + 1$$
$$= \frac{(t_{1,a} - \tau_a)(q_a \tau_a - t_{1,a} q_a + 2\Delta T_a p_a)}{2\Delta T_a p_a} \ ,$$

and again, from definition of $\Delta T_a$ and $\tau_a$ in Appendix B.1 we get

$$\frac{(t_{1,a} - \tau_a)(q_a \tau_a - t_{1,a} q_a + 2\Delta T_a p_a)}{2\Delta T_a p_a} = \frac{(t_{1,a} - t_{0,a})p_a(p_a + 1)}{2} \ . \tag{14}$$

Then, from Equations 12, 13 and 14, and recalling that $p_a + q_a = 1$, we get

$$\int_{t_{0,a}}^{t_{1,a}} \phi_a(x)dx = \frac{(t_{1,a} - t_{0,a})p_a q_a}{2} + \frac{(t_{1,a} - t_{0,a})p_a(p_a + 1)}{2}$$
$$= \frac{(t_{1,a} - t_{0,a})}{2}(p_a(1 - p_a) + p_a(1 + p_a))$$
$$= \frac{(t_{1,a} - t_{0,a})}{2}(p_a - p_a^2 + p_a + p_a^2)$$
$$= \frac{(t_{1,a} - t_{0,a})}{2}2p_a$$
$$= (t_{1,a} - t_{0,a})p_a$$
$$= \Delta T_a p_a \ ,$$

and therefore

$$\int_{t_{0,a}}^{t_{1,a}} \zeta_a(x)dx = \int_{t_{0,a}}^{t_{1,a}} \phi_a(x)dx \ ,$$

which means

$$\int_{-\infty}^{\infty} \zeta_a(x)dx = \int_{-\infty}^{\infty} \phi_a(x)dx \ .$$

□

Proofs for other curves follow the same logic.

## C OBTAINING FULL-RANK LINEAR SYSTEMS TO FIND CLOSED-FORM PIECE-WISE SOLUTIONS OF DIFFERING SMOOTHNESS

### C.1 Linear System

We search for a family of possible solutions for each group $\phi_a$, satisfying equalised odds (Definition 2.1), that adhere to the following constraints:

**continuity**

$$\phi_a(t_{0,a}) = 0 \qquad \phi_a(t_{1,a}) = 1 \, ,$$

**monotonicity**

$$\phi_a'(x) \geq 0 \, , \text{ and}$$

**preservation of probability**

$$\int_{t_{0,a}}^{t_{1,a}} \phi_a(x)dx = \int_{t_{0,a}}^{t_{1,a}} p_a dx \, .$$

By assuming that each $\phi_{a,n} = a_n + b_n x^2 + \cdots$, we can use these constraints (as well as other, more strict constraints) to find solutions to this problem of varying smoothness by solving the linear system

$$A\mathbf{x} = b \, ,$$

where $\mathbf{x} = [a_0, b_0, \ldots, a_n, c_n, \ldots]^T$. (Continuous, non-smooth solutions to this linear problem are given in Appendix B.2.)

### C.2 Closed-form Smoothness for $p_a \in [\frac{2}{5}, \frac{3}{5}]$

Here, we search for a smooth closed-form solution to the above problem. For simplicity, we assume that $t_{0,a} = 0$ and $t_{1,a} = 1$, however the solution can be generalised to arbitrary thresholds by applying shift and stretch operations.

We have the following constraints:

(1) $\psi_a(0) = 0$,
(2) $\psi_a(1) = 1$,
(3) $\psi_a'(0) = 0$,
(4) $\psi_a'(1) = 0$, and
(5) $\int_0^1 \psi_a(x)dx = \int_0^1 p_a dx = p_a$.

Having five constraints requires five coefficients, and so we assume that

$$\psi_a(x) = a_1 + b_1 x + c_1 x^2 + d_1 x^3 + g_1 x^4 \, .$$

We know that

$$\psi_a'(x) = b_1 + 2c_1 x + 3d_1 x^2 + 4g_1 x^3 \qquad \int_0^1 \psi_a(x)dx = a_1 + \frac{1}{2}b_1 + \frac{1}{3}c_1 + \frac{1}{4}d_1 + \frac{1}{5}g_1 \, ,$$

and so we have the following well-defined, full-rank linear system:

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 & 4 \\ 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix} \begin{bmatrix} a_1 \\ b_1 \\ c_1 \\ d_1 \\ g_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ p_a \end{bmatrix} \, .$$

After solving the system, we get

$$\psi_a(x) = (30p_a - 12)x^2 + (-60p_a + 28)x^3 + (30p_a - 15)x^4 .$$

If the scoring function $f(\mathbf{x}) = y$ and $y \in \mathbb{R}$, then this gives the final solution

$$\mathbb{P}\{\phi_a(y) = 1\} = \begin{cases} 0 & \text{if } y < t_{0,a} \\ \psi_a(\frac{y - t_{0,a}}{t_{1,a} - t_{0,a}}) & \text{if } t_{0,a} \le y < t_{1,a} \\ 1 & \text{if } y \ge t_{1,a} \end{cases} .$$

This function, however, is only monotonic for $p_a \in [\frac{2}{5}, \frac{3}{5}]$.

$\square$

## C.3 Piece-wise Cubic Interpolant

The closed-form, 4$^{\text{th}}$ order solution satisfies the smoothness and boundary constraints, but violates the monotonic constraint for any probability outside of the $[\frac{2}{5}, \frac{3}{5}]$ range. We can address this issue by constructing a piece-wise spline based on cubic polynomials given by

$$\psi_{a,n}(x) = a_n + b_n x + c_n x^2 + d_n x^3 .$$

However, we need to add two additional constraints to the optimisation problem: (1) an agreed meeting point and (2) an agreed derivative at the meeting point. Again, assuming that $t_{0,a} = 0$ and $t_{1,a} = 1$ – recall that we can shift and re-scale the solution later – we get:

(1) $\psi_{a,0}(0) = 0$,
(2) $\psi_{a,1}(1) = 1$,
(3) $\psi'_{a,0}(0) = 0$,
(4) $\psi'_{1,0}(1) = 0$,
(5) $\psi_{a,0}(1 - p_a) = p_a$,
(6) $\psi_{a,1}(1 - p_a) = p_a$,
(7) $\psi'_{a,0}(1 - p_a) - \psi'_{a,1}(1 - p_a) = 0$, and
(8) $\int_0^{1 - p_a} \psi_{a,0}(x)dx + \int_{1-p_a}^1 \psi_{a,1}(x)dx = \int_0^1 p_a dx = p_a$.

Since

$$\psi'_{a,n}(x) = b_n + 2c_n x + 3d_n x^2$$

and

$$\int_0^{1-p_a} \psi_{a,0}(x)dx = a_0(1 - p_a) + b_0 \frac{(1 - p_a)^2}{2} + c_0 \frac{(1 - p_a)^3}{3} + d_0 \frac{(1 - p_a)^4}{4}$$

$$\int_{1-p_a}^1 \psi_{a,1}(x)dx = a_1(1 - (1 - p_a)) + b_0 \frac{1 - (1 - p_a)^2}{2} +$$

$$c_0 \frac{1 - (1 - p_a)^3}{3} + d_0 \frac{1 - (1 - p_a)^4}{4} ,$$

we get the following linear system:

$$
A = \begin{bmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 1-p_a & (1-p_a)^2 & (1-p_a)^3 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\
0 & 0 & 0 & 0 & 0 & 1 & 2 & 3 \\
0 & 0 & 0 & 0 & 1 & 1-p_a & (1-p_a)^2 & (1-p_a)^3 \\
0 & 1 & 2(1-p_a) & 3(1-p_a)^2 & 0 & -1 & -2(1-p_a) & -3(1-p_a) \\
1-p_a & \frac{(1-p_a)^2}{2} & \frac{(1-p_a)^3}{3} & \frac{(1-p_a)^4}{4} & 1-(1-p_a) & \frac{1-(1-p_a)^2}{2} & \frac{1-(1-p_a)^3}{3} & \frac{1-(1-p_a)^4}{4}
\end{bmatrix}
$$

$$
\mathbf{x} = \begin{bmatrix} a_0 \\ b_0 \\ c_0 \\ d_0 \\ a_1 \\ b_1 \\ c_1 \\ d_1 \end{bmatrix}
\qquad
b = \begin{bmatrix} 0 \\ 0 \\ p_a \\ 1 \\ 0 \\ p_a \\ 0 \\ p_a \end{bmatrix}.
$$

Solving $A\mathbf{x} = b$ yields

$$
a_0 = 0 \qquad\qquad b_0 = 0, \qquad\qquad c_0 = \frac{3p}{p_a^2 - 2p_a + p_a} \qquad d_0 = \frac{2p_a}{(p_a - 1)^3}
$$

$$
a_1 = \frac{p_a^3 + 3p^2 - 3p + 2}{p_a^3} \qquad b_1 = \frac{6(p_a^2 - 2p + 1)}{p_a^3} \qquad c_1 = \frac{3(p_a^2 - 3p + 2)}{p_a^3} \qquad d_1 = \frac{2(p_a - 1)}{p_a^3}.
$$

We then take the solution to be

$$
\mathbb{P}\{\phi_a(y) = 1\} = \begin{cases}
0 & \text{if } y < t_{0,a} \\
\psi_{a,0}\left(\frac{y - t_{0,a}}{t_{1,a} - t_{0,a}}\right) & \text{if } t_{0,a} \le y < \tau_a \\
\psi_{a,1}\left(\frac{y - t_{0,a}}{t_{1,a} - t_{0,a}}\right) & \text{if } \tau_a \le y < t_{1,a} \\
1 & \text{if } y \ge t_{1,a}
\end{cases}.
$$

□

# D   MAXIMUM $L_{\mathcal{R}}$ FOR DIFFERENT CURVES

As discussed in Section 4.2, we know that

$$
L_{\mathcal{R}} = \max|\phi_a'(r)|.
$$

Since these curves are specified through closed-form solutions parameterised by $t_{a,y}$ and $p_a$ on a known interval $\mathcal{R}$, $L_{\mathcal{R}}$ can be found analytically for each curve. Here, we show the derivation procedure for the linear and $4^{\text{th}}$ order solutions. The other curves (cubic and quadratic) follow the same protocol.

The linear solution is defined in Equation 10. As such, we know that

$$\phi'_a(r) = \begin{cases} 0 & r \geq t_{1,a} \\ \frac{p_a}{\Delta T_a(1-p_a)} & \tau_a \leq r < t_{1,a} \\ \frac{1-p_a}{\Delta T_a p_a} & t_{0,a} \leq r < \tau_a \\ 0 & r < t_{0,a} \end{cases} .$$

Thus, the value of $L_{\mathcal{R}}$ is related to the value of $p_a$ and the distance between the thresholds with

$$\max|\phi'_a(r)| = \begin{cases} \frac{1}{\Delta T_a} & \text{if } p_a = \frac{1}{2} \\ \frac{p_a}{\Delta T_a(1-p_a)} & \text{if } p_a > \frac{1}{2} \\ \frac{1-p_a}{\Delta T_a p_a} & \text{if } p_a < \frac{1}{2} \end{cases} .$$

The 4$^{\text{th}}$ order is define in Equation 5, and takes the form

$$\phi_a(r) = \begin{cases} 1 & r \geq t_{1,a} \\ \psi_a\left(\frac{r-t_{0,a}}{t_{1,a}-t_{0,a}}\right) & t_{0,a} \leq r < t_{1,a} \\ 0 & r < t_{0,a} \end{cases} ,$$

where

$$\psi_a(x) = (30p_a - 12)x^2 + (-60p_a + 28)x^3 + (30p_a - 15)x^4 . \tag{15}$$

Since the definition of $\phi_a(r)$ is always constrained such that it is monotonic and $\phi_a(t_{y,a}) = y$, the maximum derivative always occurs at the point of inflexion, or

$$\psi''_a\left(\frac{r - t_{0,a}}{t_{1,a} - t_{0,a}}\right) = 0 .$$

From Equation 15:

$$\psi''_a(x) = 12(30p_a - 15)x^2 + 6(-60p_a + 28)x + 2(30p_a - 12) ,$$

and so from the quadratic formula for $\psi''_a(x) = 0$ we get

$$x^\star = \frac{-(6(-60p_a + 28)) - \sqrt{(6(-60p_a + 28))^2 - 4(12(30p_a - 15))(2(30p_a - 12))}}{2(12(30p_a - 15))}$$

$$= -\frac{7 - 15p_a + \sqrt{75p_a^2 - 75p_a + 19}}{30p_a - 15} ,$$

where we ignore the second root, as it takes $x^\star$ out of range. We can use this to find $\max|\phi'_a(x)| = \phi'_a(x^\star)$, which assumes $t_{0,a}$ and $t_{1,a}$ are fixed at 0 and 1 respectively, and apply a re-scaling such that

$$\max|\phi'_a(r)| = \frac{|\phi'_a(x^\star)|}{t_{1,a} - t_{0,a}} .$$

This method works for all valid values of $p_a$, except when $p_a = \frac{1}{2}$ where $\psi_a(x)$ is reduced to a cubic equation, and thus the point of inflexion is perfectly between the two thresholds:

$$\max|\phi'_a(r)| = \phi'_a\left(t_{0,a} + \frac{1}{2}(t_{1,a} - t_{0,a})\right) ,$$

in which case

$$\mathcal{L}_{\mathcal{R}} = \max |\phi'_a(r)| \, .$$
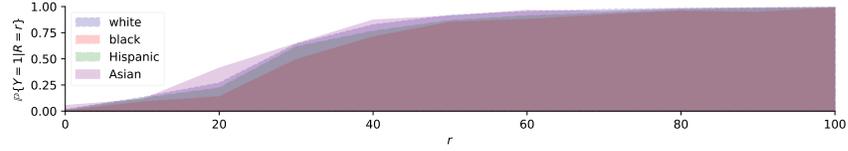
## E PROBABILITY FUNCTIONS



Fig. 8. Probability function for the CreditRisk data set. An increase in credit score $r$ (x-axis) generally leads to an increase in the probability of an individual not defaulting on a loan (y-axis) in the last 90 days ($Y = 1$).
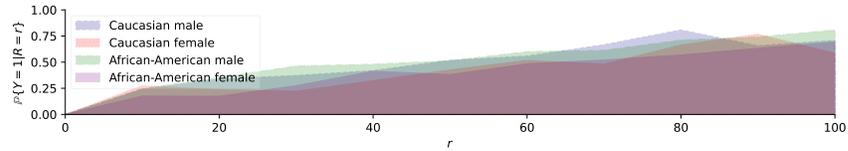


Fig. 9. Probability function for the COMPAS data set. An increase in score $r$ (x-axis) generally leads to an increase in the probability of committing an offence (y-axis) in a two-year time window ($Y = 1$).

## F OPTIMAL THRESHOLDS FOR FAIRNESS

| | white | | | black | | | Hispanic | | | Asian | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $t_{0,a}$ | $t_{1,a}$ | $p_a$ | $t_{0,a}$ | $t_{1,a}$ | $p_a$ | $t_{0,a}$ | $t_{1,a}$ | $p_a$ | $t_{0,a}$ | $t_{1,a}$ | $p_a$ |
| fixed | 25.0 | 50.0 | 0.500 | 23.5 | 85.0 | 0.972 | 30.0 | 30.0 | 0.000 | 35.5 | 49.0 | 0.728 |
| linear | 9.5 | 52.0 | 0.348 | 20.5 | 42.5 | 0.830 | 30.0 | 30.0 | 0.000 | 33.5 | 61.5 | 0.806 |
| quadratic | 11.0 | 78.5 | 0.610 | 20.0 | 82.0 | 0.928 | 30.0 | 30.0 | 0.000 | 24.0 | 49.5 | 0.406 |
| cubic | 1.5 | 49.5 | 0.256 | 20.5 | 44.5 | 0.844 | 30.0 | 30.0 | 0.000 | 34.0 | 70.0 | 0.864 |
| $4^{\text{th}}$ order | 7.5 | 63.5 | 0.468 | 14.0 | 32.0 | 0.426 | 30.0 | 30.0 | 0.000 | 28.0 | 52.5 | 0.546 |

Table 3. Thresholds and probabilities for each curve across all classes of the CreditRisk data set. See Figure 4 for visualisation.

| | Caucasian male | | | Caucasian female | | | African-American male | | | African-American female | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $t_{0,a}$ | $t_{1,a}$ | $p_a$ | $t_{0,a}$ | $t_{1,a}$ | $p_a$ | $t_{0,a}$ | $t_{1,a}$ | $p_a$ | $t_{0,a}$ | $t_{1,a}$ | $p_a$ |
| fixed | 24.0 | 41.0 | 0.116 | 7.0 | 37.0 | 0.048 | 48.0 | 48.0 | 0.000 | 35.0 | 58.0 | 0.930 |
| linear | 20.0 | 43.0 | 0.194 | 27.0 | 42.0 | 0.478 | 48.0 | 48.0 | 0.000 | 23.0 | 42.0 | 0.326 |
| quadratic | 18.0 | 46.0 | 0.266 | 26.0 | 47.0 | 0.576 | 48.0 | 48.0 | 0.000 | 12.0 | 43.0 | 0.232 |
| cubic | 34.0 | 54.0 | 0.772 | 33.0 | 91.0 | 0.952 | 30.0 | 30.0 | 0.000 | 33.0 | 77.0 | 0.926 |
| $4^{th}$ order | 25.0 | 48.0 | 0.412 | 26.0 | 46.0 | 0.557 | 48.0 | 48.0 | 0.000 | 26.0 | 48.0 | 0.550 |

Table 4. Thresholds and probabilities for each curve across all classes of the COMPAS data set. See Figure 6 for visualisation.