

---

# FairCanary: Rapid Continuous Explainable Fairness

---

Avijit Ghosh\*

Northeastern University, Fiddler Labs<sup>†</sup>  
avijit@ccs.neu.edu

Aalok Shanbhag\*

Fiddler Labs  
aalok@fiddler.ai

## Abstract

Machine Learning (ML) models are being used in all facets of today's society to make high stake decisions like bail granting or credit lending, with very minimal regulations. Such systems are extremely vulnerable to both propagating and amplifying social biases, and have therefore been subject to growing research interest. One of the main issues with conventional fairness metrics is their narrow definitions which hide the complete extent of the bias by focusing primarily on positive and/or negative outcomes, whilst not paying attention to the overall distributional shape. Moreover, these metrics are often contradictory to each other, are severely restrained by the contextual and legal landscape of the problem, have technical constraints like poor support for continuous outputs, the requirement of class labels, and are not explainable.

In this paper, we present Quantile Demographic Drift, which addresses the shortcomings mentioned above. This metric can also be used to measure intra-group privilege. It is easily interpretable via existing attribution techniques, and also extends naturally to individual fairness via the principle of like-for-like comparison. We make this new fairness score the basis of a new system that is designed to detect bias in production ML models without the need for labels. We call the system FairCanary because of its capability to detect bias in a live deployed model and narrow down the alert to the responsible set of features, like the proverbial canary in a coal mine.

## 1 Introduction

Algorithmic Fairness research has seen a rising growth in academic interest in recent times. The application of machine learning for sensitive tasks like criminal sentencing [3], job marketplaces [17], commercial facial recognition algorithms [15], and health systems [54] has naturally compelled researchers to unearth the disparities produced by these algorithms. Several competing definitions of fairness have been discussed in the literature, which according to Corbett-Davies et Al.[19], fall under three general classes of work: (1) *Anti-classification*, where protected features or their proxies are not used to make the decision, (2) *Classification Parity*, where common measures of model predictive performance are equal across protected groups and (3) *Calibration*, where the outcomes, conditional on priors, are independent of protected features. They also dissect these ideas and metrics, claiming that they have "deep statistical limitations"[19], with several metrics at odds with one another. Furthermore, a large number of conventional metrics depend on different combinations of values from the confusion matrix of a binary classifier, which carries with it a host of problems including non-generalizability due to choice of threshold, and the tension between utility and representativeness in the classification output. In addition, models which have a continuous output, instead of discrete output classes, or have sensitive protected attributes that happen to be continuous values, these simple ratio based metrics simply fail to contextualize these nuances.

---

\*Joint first authors.

<sup>†</sup>Work done as an intern at Fiddler Labs.

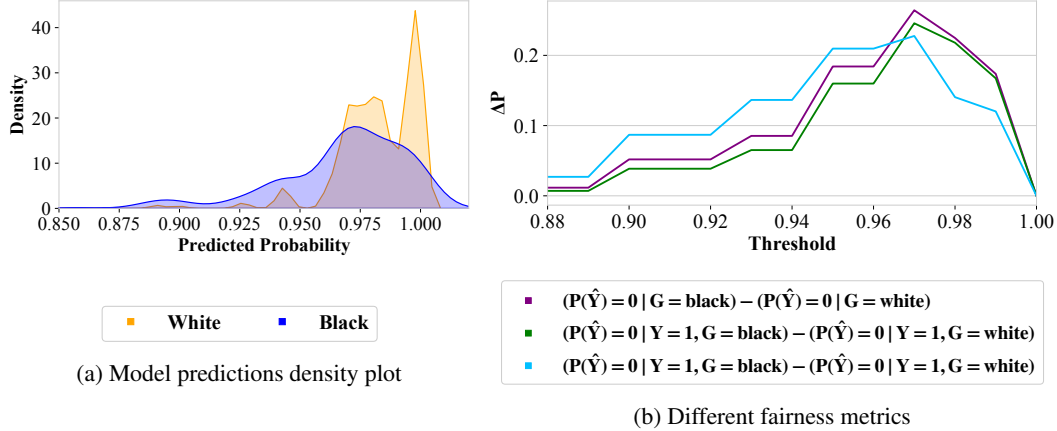


Figure 1: A plot, similar to [48], showing three threshold-based fairness metrics (statistical parity, equalized odds, equal opportunity) on the LSAC admissions dataset [69]. The choice of threshold greatly impacts the fairness metric and hides the complete context of the bias.

Threshold-independent fairness metrics hold the clue towards beginning to solve this dual conundrum. In the case of continuous outputs, past work has suggested measuring the full spectrum of prediction probability distributions, naturally also operating independent of decision labels, for instance KL Divergence, JS Divergence, and Wasserstein distance. In particular, Wasserstein distance has been explored as a fairness metric because of its optimal transport characteristics and its differentiable properties that lead itself to SHAP-like explanations. [48]. On the other hand, in the case of continuous features, the Hirschfeld Gebelein-Renyi (HGR) Maximum Correlation Coefficient and Decision Tree based methods [58] have been proposed to measure the independence of continuous features against a continuous output.

In a production scenario with continuously arriving data points, even a trained "fair" classification model might become unfair over time with the presence of new, unseen data. Without access to the ground truth labels of this new data, measuring and mitigating unfairness in such a continuous learning scenario has been scarcely studied in FairML literature. In addition, another desirable property of such a monitoring system would be explainability - if a model begins to be measured as unfair, a truly robust monitoring system would also point out the features that are the most responsible for the change in the fairness metric, and also the regions in the stream of new data where this unfairness occurs.

We aim to answer the following questions for a machine learning model in production: (1) *Which biases can be measured, given labels may or may not be present*, (2) *How to measure these biases*, (3) *How to explain these biases, in terms of the input features of the model*, and (4) *How to mitigate these biases*. We answer the first question by identifying population bias, algorithmic bias and individual bias [46] as the three relevant biases for a machine learning model in production, that can be measured in real time. We also identify Intra-group privilege bias as an additional useful metric, because of the possibility of stark disparities within a protected group itself that might not be detected by overall group fairness metrics. We then introduce a novel quantile based measure of distributional difference to measure algorithmic bias. Incidentally, there already exists literature to mitigate bias which is similar in end result to quantile norming Nanda et al. [50]. We use the quantile alignment that we gain naturally to measure individual bias, and also for mitigation. Algorithmic and individual bias is then explained using standard attribution techniques that satisfy the efficiency axiom [44, 66].

## 2 Background

### 2.1 Algorithmic Fairness

The use of algorithms to aid critical decision making processes in the government and the industry has attracted commensurate scrutiny from academia, lawmakers and social justice workers in recent times [4, 7, 71], because ML systems trained on a snapshot of the society has the unintended

consequences of learning, propagating and amplifying historical social biases and power dynamics [5, 56]. The current research landscape consists of both ML explanation methods and fairness metrics to try and uncover the problems of trained models [8, 30, 45, 59, 68], and fairness aware ML algorithms, for instance classification [31, 34, 37, 47], regression [2, 9], causal inference [43, 49], word embeddings [13, 14] and ranking [16, 64, 72].

## 2.2 Concept Drift and Fairness Monitoring

Any deployed machine learning model cannot guarantee consistent performance over time, if the underlying data changes stochastically. This phenomenon, called concept drift, has been well studied in the literature, from sudden drifts [53] to gradual drifts [65]. Furthermore, such drifts may be differentiated between a true shift in the relationship between the underlying variables, or may be attributed to a sampling issue [60]. However, most studies looking at drift look at traditional performance metrics like accuracy or recall. In a scenario where a deployed machine learning model is making sensitive decisions, we posit that it is of considerable interest to systematically measure and analyse the drift in the fairness performance of the model.

While most popular methods for detecting concept drift [10, 27] assume that the labels for the predicted variable are immediately available, this might not be feasible or even accurate if an actual concept drift has taken place, rendering the final labels unreliable. A class of prior work [23, 29, 57, 73] measures the drift of prediction distributions as a proxy for concept drift. This is class of work is the most aligned with our proposed method.

The need for measuring and mitigating unfairness in deployed systems is also being increasing backed by legislative action around the world. The recently proposed Artificial Intelligence Act by the European Commission [18] stresses on this aspect - *"The provider should establish a sound quality management system, ensure the accomplishment of the required conformity assessment procedure, draw up the relevant documentation and establish a robust post-market monitoring system."*, and more specifically to fairness, *"the providers should be able to process also special categories of personal data, as a matter of substantial public interest, in order to ensure the bias monitoring, detection and correction in relation to high-risk AI systems."*. Similar legislation has been proposed in New Zealand [39], Canada [55], the US [1] and the UK [25]. All these point to the requirement of a robust and explainable bias monitoring system for deployed machine learning models.

## 2.3 Interpretability of model fairness

Interpretation of fairness metrics in terms of the input features to the model has not been studied extensively so far. Miroshnikov et al. [48] proposed explaining the Wasserstein-1 distance using a Shapley value formulation. One of the issues with this approach is that the attributions need to be computed  ${}^nC_2$  times if there are  $n$  protected groups. Also, it can be computationally challenging to compute Shapley values over large samples. [63] proposed explaining differentiable distance metrics using Integrated Gradients, but this needs the model to be differentiable as well.

Explaining the standard fairness metrics that use ground truth labels using Shapley values is possible by making the assumption that the perturbed values retain the original output label. But this could be misleading because the perturbations change the nature of the instance, and can even create Out-of-Distribution (OOD) points [40].

Interpretation of fairness metrics in terms of the constituent features is a very useful tool to analyze the causes of unfairness of the model. Note that causes here refers to causality in terms of the model, not in terms of the actual, population level causal relationship between the inputs and outputs.

## 3 Shortcomings of existing fairness metrics

**Subjectivity of Algorithmic fairness** The field of algorithmic fairness heavily derives from interdisciplinary schools of thoughts and as such, is the focus of philosophical debates around the task of objectively defining something as subjective as fairness. Binns [11] draw definitions of machine learning fairness from political philosophy, while Hutchinson and Mitchell [35] trace the history of fairness definitions in education and machine learning. Sambasivan et al. [61] point out the non-compatibility of US-centric fair machine learning research in an eastern society like

| Metric/Framework                    | Related terms                                                                        | Continuous groups | Continuous outputs | Interpretability |
|-------------------------------------|--------------------------------------------------------------------------------------|-------------------|--------------------|------------------|
| Demographic parity [24]             | mean difference, demographic parity, disparate treatment, group discrimination score | ✗                 | ✗                  | ✗                |
| Conditional statistical parity [20] | statistical parity, conditional procedure accuracy, disparate treatment              | ✗                 | ✗                  | ✗                |
| Equalized odds [33]                 | equalized odds, false positive/negative parity, disparate treatment                  | ✗                 | ✗                  | ✗                |
| Equal opportunity [33]              | equality of opportunity, individual fairness, disparate treatment                    | ✗                 | ✗                  | ✗                |
| Counterfactual fairness [12, 41]    | counterfactual fairness, disparate treatment, fliptest                               | ✗                 | ✗                  | ✗                |
| Statistical independence [32]       | HGR coefficient, independence                                                        | ✓                 | ✓                  | ✗                |
| Distributional difference [48]      | KL divergence, JS Divergence, Wasserstein distance                                   | ✓                 | ✓                  | ✓                |

Table 1: A summary showing the pros and cons of different classes of conventional fairness metrics in the literature. The metric families are largely inspired by Mehrabi et al. [46]. The related terminology column is from Das et al. [21].

India because of the differing notions of justice (distributive justice vs social justice). Also different legal systems around the world mandate different fairness goals – US Law mandates the absence of Disparate Impact [70] in hiring, while Indian Law provides for fixed quotas in areas of opportunity like higher education and government employment [22].

**Statistical, Temporal and Interpretability Limitations** Conventional fairness metrics have impossibility results [52], meaning that they are often at odds with each other. Previous work [19, 38, 48] point out that it is impossible to satisfy both *Classification parity* and *Calibration* metrics at the same time, except for very specific conditions, and therefore context becomes key when picking a metric [6, 62]. The statistical limitations extend to group membership limitations - most conventional metrics require groups and subgroups to be discrete variables and cannot work with continuous variables [30], and "confusion matrix based metrics" [52] additionally do not support continuous outputs (which is often the case in problems like regression and recommendation), causing measurement to be severely limited with ad-hoc thresholds that causes interpretation to wildly differ (Figure 1). Also, most conventional fairness metrics ignore stochasticity – a temporal analysis in [42] showed how the changing of fairness metrics over time, due to data drift, concept drift or otherwise, could actually harm sensitive groups, especially when redressal is based on a measurement from a fixed point in time. Additionally, because these measurements are often in-situ and lacking context, it is challenging to interpret the underlying causes of bias from the metric itself. Only one work [48], from recent literature, shows how one can break down a Wasserstein distance based fairness metric into Shap-like feature attributions. Table 1 shows an overview of the terminology and limitations of different classes of fairness metrics in the literature.

## 4 Distributional Distance as a Fairness Metric

There are several reasons for measuring the fairness of machine learning model by comparing segments of its prediction distribution to other segments, or to the baseline population distribution of the segment. The primary reason is to ensure that we measure bias across the board, as opposed to focusing on the positive selected instances, or on group level approximations. As groups get smaller, they start revealing more information about intra group disparities that would have otherwise been lost due to aggregation [30]. The measure should hence reflect the bias across the sample, as we divide it into smaller and smaller buckets, right down to individual instances. This helps remove aggregation bias from the bias measurement itself.

#### 4.1 Desired properties of a metric to measure bias

We now discuss certain desirable properties of a distributional fairness metric that fits our stated objectives:

1. The metric should be in the units of the model’s prediction. The utility of this is especially evident when dealing with continuous output models.
2. It should take the value zero only if the observed prediction distributions being compared are exactly the same.
3. It should be continuous with respect to the change in the geometry of the distribution [48].
4. It should be non-invariant with respect to monotone transformations of the distributions [48]. For example, given two samples of points  $S_1$  and  $S_2$ , if we multiply the value of each point in the samples by a constant  $k$ , the distance between the modified samples should now depend on the  $k$ . JSD for example does not satisfy this property.

The metric should also be bias transforming as described in [67] i.e. should not be satisfied by a model that preserves the biases present in the data.

### 5 Quantile Demographic Disparity

We now describe an alternative metric called Quantile Demographic Disparity (QDD), which is a function of the quantile bin that a prediction event lies in. QDD is designed to work for continuous outputs and can be customized to provide sliced views down to an individual level if necessary. As we explain later, QDD is inherently interpretable, and easy for stakeholders, (who may not be ML practioners), to understand.

For the two groups  $G_1$  and  $G_2$ , let the two distributional samples of model scores be  $S_1$  and  $S_2$ . We divide the samples into  $B$  bins, of equal size  $N_1$  and  $N_2$  respectively. This is equivalent to segmenting by quantiles. For example, if there are 10 bins, we are essentially bucketing individuals between the 0th and 10th percentile, 10th-20th percentile and so on.

QDD for bin  $b$  is defined as:

$$QDD_b = \mathbb{E}_{G_{1,b}}[S_1] - \mathbb{E}_{G_{2,b}}[S_2] \quad (1)$$

This can be approximated as:

$$QDD_b = \frac{1}{N_1} \sum_{n=1}^{N_1} S_{1,n} - \frac{1}{N_2} \sum_{n=1}^{N_2} S_{2,n} \quad (2)$$

The QDD, when conditioned on certain attributes  $C$ , becomes the Conditional Quantile Demographic Disparity.

We introduce several perspectives that can be obtained via QDD calculations.

**Intra-Group bias** The maximum difference in the extent of bias within a group, when comparing between bins. This helps us disentangle intra-group aggregation bias. This is defined as the maximum QDD across the  $b$  bins.

**Disparity with base rate** The difference in the extent of bias when compared to the bias of the population. This can be defined as the difference of the QDD when calculating it over the production data and the QDD of the training data.

**Individual Fairness via Alignment** The QDD is defined between two protected class groups, over a given number of bins, which determines the resolution of the metric. If the number of bins is equal to the number of instances in the sample, we are left with comparisons between individuals at the same rank or percentile comparisons. This is equivalent to the concept of alignment proposed in Shanbhag et al. [63]. This gives us a clean way to obtain Individual fairness insights. The counterfactual

example here is the same ranked counterpart in the opposite group. This method doesn't require us to compute complex counterfactuals, which could have their own biases and errors. The principle we use to justify this is as follows. If there is no bias between the two groups, and we have a large enough sample of both, then in the prediction space, which is a task specific dimensionality reduction, the distance between individuals of the same rank should be zero.

**Local Quantile Disparity Attribution** The local quantile disparity attribution for a feature  $f$  explains the local quantile disparity in terms of the feature  $f$ , using an attribution method  $A$  that satisfies the axiom of efficiency.

**Definition** The QDD Attribution for feature  $f$ , for prediction sample  $S_1$  over  $S_2$  in bin  $b$ ,  $QDDA_{b,A,f}$  is a measure of the change in QDD in bin  $b$  that can be attributed to feature  $f$  using attribution method  $A$  that satisfies the efficiency axiom.

$$QDDA_{b,A,f} = \frac{1}{N_t} \sum_{n=1}^{N_t} attr_{n,A,S_1} - \frac{1}{N_r} \sum_{n=1}^{N_r} attr_{n,A,S_2} \quad (3)$$

where  $attr_{n,A,S_1}$  refers to the attribution of the  $n$ th data point to feature  $f$  for a prediction from bin  $b$  of distribution  $S_1$  using attribution method  $A$ .

Given that the attribution method  $A$  satisfies the efficiency axiom,  $QDD_b = \sum_{f=1}^F QDDA_{b,A,f}$ . Proof is provided in the Appendix.

Explaining the disparity in this manner enables a single attribution to be used for multiple explanations across groups, as compared to Shapley values over a particular metric, which need to be re-calculated for every grouping. Our explanation technique therefore is much less calculation intensive than previous techniques like [48], since it requires the calculation of attributions only once.

## 6 FairCanary

We now describe FairCanary, a system to address the questions posed earlier. We will not go into the technical details of calculating metrics for streaming data, as these calculations are fairly trivial.

### 6.1 Measurement

First, we will identify the biases, and suggest appropriate metrics for measurement.

1. Population bias - Measured via difference in representation in the real world, via the dataset and in production. This assumes that the representation in the dataset is fair. Here we store counts for each intersectional sub-group.
2. Algorithmic/Historical bias - Measure QDD over  $b$  bins for each intersectional sub-group
3. Intra-group privilege - for each sub-group, we calculate the maximum QDD difference between two bins. This is relevant from the perspective of recourse. Say a model satisfies the accuracy requirements for fairness for a given threshold, but is unfair for individuals in a certain group by giving them depressed scores. Then this has implications for the recourse that is suggested, the individual may be suggested a recourse that is too harsh for them. This would also be a kind of disparate treatment.
4. Individual bias - for a case flagged for view, either by the recipient of the model's decision, or by a human-in-the-loop, we compare against the instances of same rank in the dominant sub-group.

As part of the measurement framework, we propose setting up alerts around each metric. This would help to determine when mitigation, or human-in-the-loop intervention is merited.

### 6.2 Attribution

For any model, we can explain the prediction output via Shapley value based methods [44], and if the model is differentiable, we can use Integrated Gradients (IG) [66]. The reason for choosing

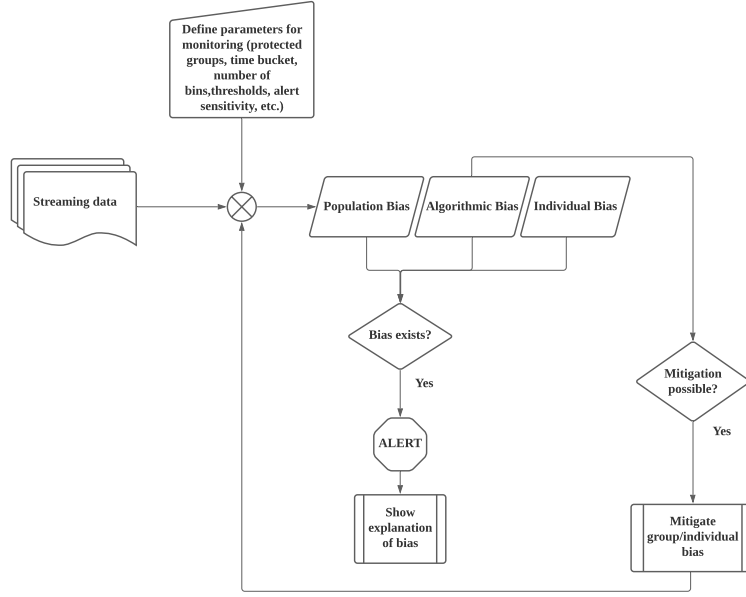


Figure 2: A schematic showing how FairCanary works in a production system.

these methods is that satisfy certain desirable axioms, one of which is efficiency, which helps with a precise accounting of the bias. It must be noted that there are other explanation methods which satisfy efficiency, but we will not explore those in detail.

It is important to know the input features that are causing the distributional difference between two groups, to help understand either correlations with protected features, or a differential in distributional change for certain features with respect to the two groups. Using Shapley values or IG, we can explain a variety of model types such as text or image classification, or tabular regression etc, and perform further operations such as feature groupings, given that the attributions from these methods can be summed up.

### 6.3 Mitigation

Mitigation is a key part of monitoring bias, enabling corrective action to be taken, if valid reasons exist for doing so. The advantages of post-processing mitigation, as opposed to pre-training debiasing are discussed in [28], Care must be taken to ensure that the mitigation is undertaken only after a holistic understanding of it’s consequences. [20] demonstrates several cases where mitigation may cause more harm to the individual or to the community of the particular sub-group. The mitigation we suggest is simply replacing the score with the score of the corresponding rank in the dominant group. This is similar to the mitigation proposed in [36, 51]. The justification for quantile norming lies in that if bias is known to exist and is the only rational explanation for disparity, and can be assumed to be equal within the protected group, then measuring the rank is a reasonable approach. Distribution norming as a mitigation method, however, might potentially be illegal in certain jurisdictions, specifically the US [35], but not in other places [18]. The legal implications of this method thus need to be revisited based on the jurisdiction where this is being used, and compared to other affirmative action methods that exist or are mandated by law [61] that could have a similar corrective effect.

To summarize, FairCanary works as follows:

1. Define all the intersectional groups of interest
2. Create the base rate statistics for all groups - representation of each group and quantile bin means for the chosen number of bins
3. Maintain counts over time for each group, for a chosen granularity of time bins
4. Metric calculation for incoming traffic for chosen time bin granularity

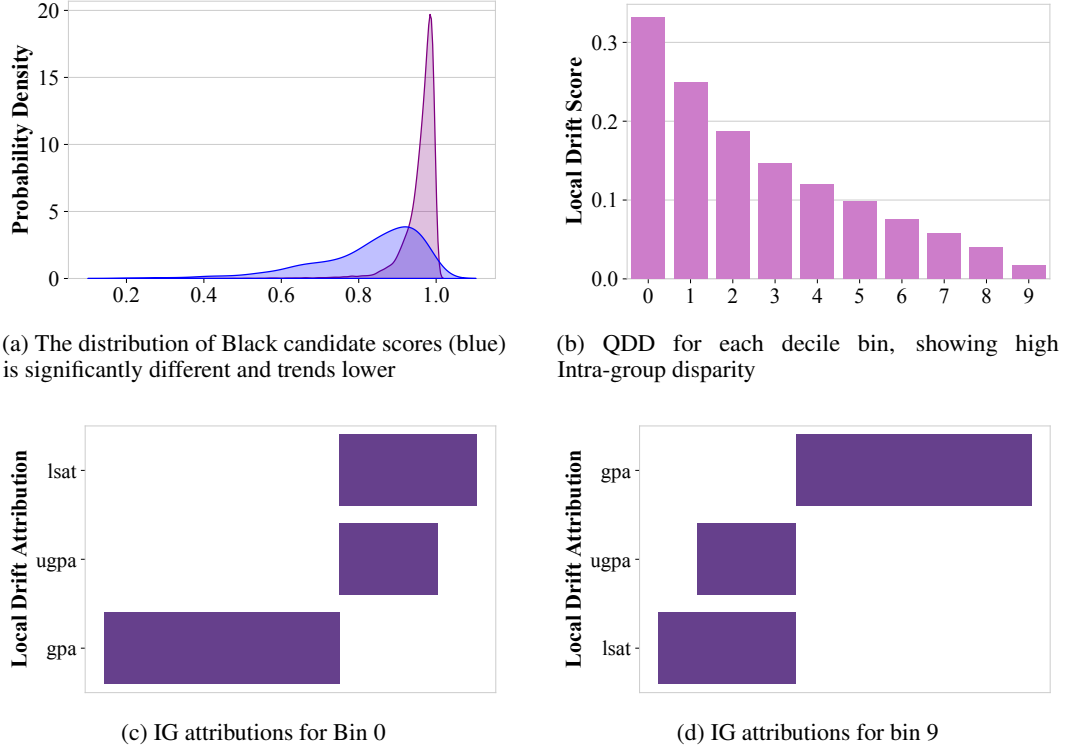


Figure 3: Prediction distributions and attributions for the LSAC dataset [69] when comparing White and Black candidates

5. Alerting across intersectional groups - if any group performs below a given threshold, alert
6. Explainability via attributions
7. Fixing the disparity - in cases where intervention is justified, use the bin prediction of the dominant class, instead of the actual prediction.

## 7 Case Study

We demonstrate the utility of this framework with a couple of brief discussions using two datasets. One is the Luxembourg dataset [74] constructed using European Social Survey data. Each instance is an individual. The labels indicate high or low internet usage, which has a continuous output. The task for this model is to predict the Internet usage of an individual. The second dataset is the LSAC dataset [69] for Law School admissions. More detailed analysis is provided in the Appendix.

For the Luxembourg dataset, we train the model on data from 2018 and compare the drift for the prediction distributions for Gender 1 [Fig 3a]. In Figure 3b, We can see there is a significant difference in the change between bins, possibly due to few very high volume users. The explanations obtained using IG in Figures 3c and 3d, help us understand this drift. The metric is signed, meaning it also tells us about the directionality of the change. Note that most existing fairness metrics don't support continuous outputs. This can be problematic, as continuous output models can have fairness implications, for example a model which suggests the interest rate to be offered to a loan applicant.

For the LSAC data, we train a model to predict the candidate passing the bar exam. We compare the QDD between a sample of black and white students. Note the intra-group disparity. If we focus simply on the aggregate numbers, it may happen in cases that disparate treatment is hidden because of the lack of bias for the privileged group within a class.



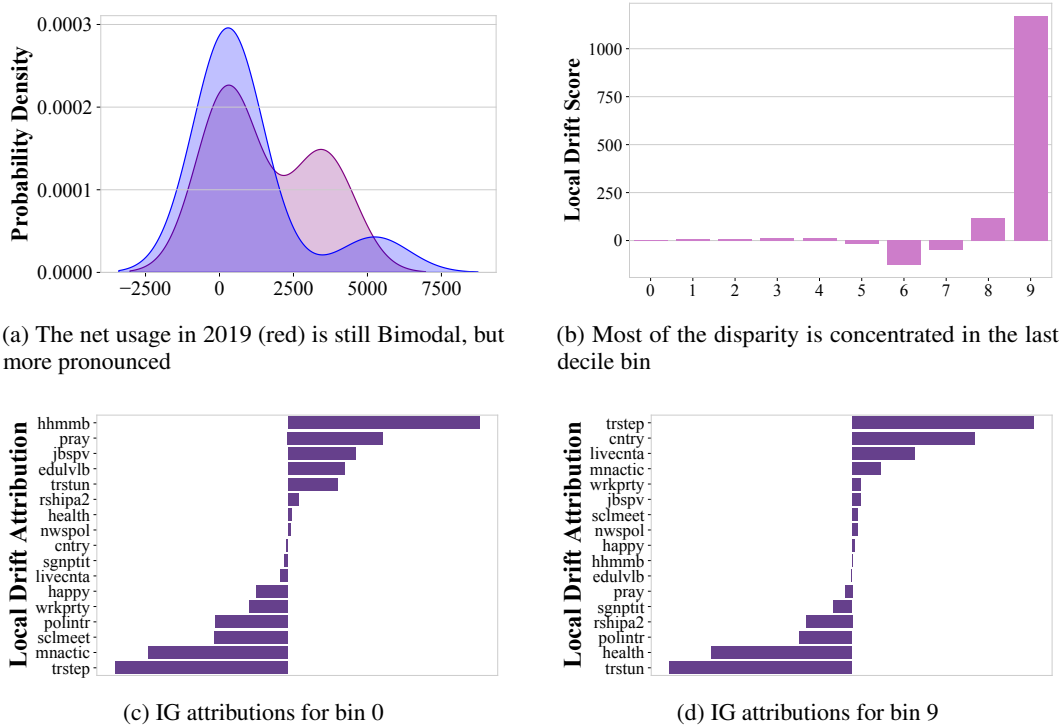


Figure 4: Prediction distributions and attributions for the Luxembourg dataset [74] when comparing Gender 1 net usage in 2018 and 2019

## 8 Discussion

**Conclusion** In this work, we present FairCanary, a comprehensive end to end framework to monitor machine learning models in production, without the need for prediction labels and threshold values. To enable this, we introduce a threshold independent metric called Quantile Demographic Disparity, and we show how QDD can be broken down into constituent feature level explanations with an approximation that is computationally much cheaper than earlier methods. We verify that FairCanary works reasonably with the help of two case studies.

**Limitations and Future Work** While threshold independence is the strength of QDD, it is also a potential weakness: Without ground truth labels, calculated disparities are, at the end of the day, best case approximations of the discrimination that actually takes place in society. We therefore do not advocate for the elimination of classical metrics that require ground truth labels and thresholds, but instead propose using them in conjunction with FairCanary to get the picture of real life harms in a context dependent manner [26]. Additionally, QDD/FairCanary is also not completely automatic, as showing in Figure 2, there are still manual parameters that need to be set, like number of bins or alert sensitivity. Training and pre-deployment stages of a model could be brought under a broadened version of this framework. Also, a fairness monitoring solution should consider providing actionable recourse explanations to the end-user via a suitable interface.

## 9 Broader Impact

We hope this framework will be adopted by companies and institutions that rely on large deployed machine learning models who are currently testing retroactively for fairness over discrete snapshots of their datasets. The Faircanary framework will let practitioners have real time visibility into their pipelines who can then proactively apply fairness interventions, which will in turn bring more equity and justice to the stakeholders impacted by large scale deployed models. This tool could also serve as a framework for regulating agencies who can get alerted of potential rogue models in real time.

## References

- [1] 116th Congress (2019-2020). H.r.2231 - algorithmic accountability act of 2019. <https://www.congress.gov/bill/116th-congress/house-bill/2231>.
- [2] Alekh Agarwal, Miroslav Dudík, and Zhiwei Steven Wu. Fair regression: Quantitative definitions and reduction-based algorithms. *arXiv preprint arXiv:1905.12843*, 2019.
- [3] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias. *propublica*. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>, 2016.
- [4] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias: There’s software used across the country to predict future criminals and it’s biased against blacks. 2016. URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>, 2019.
- [5] Solon Barocas and Andrew D. Selbst. Big data’s disparate impact. *104 California Law Review*, 671, 2016.
- [6] Solon Barocas, Anhong Guo, Ece Kamar, Jacquelyn Krones, Meredith Ringel Morris, Jennifer Wortman Vaughan, Duncan Wadsworth, and Hanna Wallach. Designing disaggregated evaluations of ai systems: Choices, considerations, and tradeoffs. *arXiv preprint arXiv:2103.06076*, 2021.
- [7] Thorsten Beck, Patrick Behr, and Andreas Madestam. Sex and credit: Is there a gender bias in lending? *Journal of Banking and Finance*, 87, 2018.
- [8] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, et al. Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943*, 2018.
- [9] Richard Berk, Hoda Heidari, Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. A convex framework for fair regression. *arXiv preprint arXiv:1706.02409*, 2017.
- [10] Albert Bifet and Ricard Gavaldà. Learning from time-changing data with adaptive windowing. In *Proceedings of the 2007 SIAM international conference on data mining*, pages 443–448. SIAM, 2007.
- [11] Reuben Binns. Fairness in machine learning: Lessons from political philosophy. In *Conference on Fairness, Accountability and Transparency*, pages 149–159. PMLR, 2018.
- [12] Emily Black, Samuel Yeom, and Matt Fredrikson. Fliptest: fairness testing via optimal transport. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 111–121, 2020.
- [13] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357, 2016.
- [14] Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard Zemel. Understanding the origins of bias in word embeddings. In *International Conference on Machine Learning*, pages 803–811, 2019.
- [15] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91, 2018.
- [16] L Elisa Celis, Damian Straszak, and Nisheeth K Vishnoi. Ranking with fairness constraints. In *45th International Colloquium on Automata, Languages, and Programming (ICALP 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- [17] Le Chen, Ruijun Ma, Anikó Hannák, and Christo Wilson. Investigating the impact of gender on rank in resume search engines. In *Proceedings of the 2018 chi conference on human factors in computing systems*, pages 1–14, 2018.
- [18] European Commission. Proposal for a regulation laying down harmonised rules on artificial intelligence (artificial intelligence act). <https://digital-strategy.ec.europa.eu/en/library/proposal-regulation-laying-down-harmonised-rules-artificial-intelligence-artificial-intelligence>.
- [19] Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*, 2018.

- [20] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*, pages 797–806, 2017.
- [21] Sanjiv Das, Michele Donini, Jason Gelman, Kevin Haas, Mila Hardt, Jared Katzman, Krishnaram Kenthapadi, Pedro Larroy, Pinar Yilmaz, and Bilal Zafar. Fairness measures for machine learning in finance.
- [22] Frank De Zwart. The logic of affirmative action: Caste, class and quotas in india. *Acta Sociologica*, 43(3): 235–249, 2000.
- [23] Denis Moreira dos Reis, Peter Flach, Stan Matwin, and Gustavo Batista. Fast unsupervised online drift detection using incremental kolmogorov-smirnov test. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1545–1554, 2016.
- [24] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226, 2012.
- [25] UK Office for Artificial Intelligence. Ethics, transparency and accountability framework for automated decision-making. <https://www.gov.uk/government/publications/ethics-transparency-and-accountability-framework-for-automated-decision-making>.
- [26] Center for Data Science and Public Policy. Aequitas: Fairness tree. <http://www.datasciencepublicpolicy.org/projects/aequitas/>.
- [27] Joao Gama, Raquel Sebastiao, and Pedro Pereira Rodrigues. On evaluating stream learning algorithms. *Machine learning*, 90(3):317–346, 2013.
- [28] Sahin Cem Geyik, Stuart Ambler, and Krishnaram Kenthapadi. Fairness-aware ranking in search & recommendation systems with application to linkedin talent search. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2221–2231, 2019.
- [29] Sindhu Ghanta, Sriram Subramanian, Lior Khremosh, Swaminathan Sundararaman, Harshil Shah, Yakov Goldberg, Drew S. Roselli, and Nisha Talagala. ML health: Fitness tracking for production models. *CoRR*, abs/1902.02808, 2019. URL <http://arxiv.org/abs/1902.02808>.
- [30] Avijit Ghosh, Lea Genuit, and Mary Reagan. Characterizing intersectional group fairness with worst-case comparisons. *arXiv preprint arXiv:2101.01673*, 2021.
- [31] Naman Goel, Mohammad Yaghini, and Boi Faltings. Non-discriminatory machine learning through convex fairness criteria. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [32] Vincent Grari, Boris Ruf, Sylvain Lamprier, and Marcin Detyniecki. Fairness-aware neural rényi minimization for continuous features. *arXiv preprint arXiv:1911.04929*, 2019.
- [33] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. *arXiv preprint arXiv:1610.02413*, 2016.
- [34] Lingxiao Huang and Nisheeth K Vishnoi. Stable and fair classification. *arXiv preprint arXiv:1902.07823*, 2019.
- [35] Ben Hutchinson and Margaret Mitchell. 50 years of test (un) fairness: Lessons for machine learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 49–58, 2019.
- [36] Ray Jiang, Aldo Pacchiano, Tom Stepleton, Heinrich Jiang, and Silvia Chiappa. Wasserstein fair classification. In *Uncertainty in Artificial Intelligence*, pages 862–872. PMLR, 2020.
- [37] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-aware classifier with prejudice remover regularizer. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 35–50. Springer, 2012.
- [38] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*, 2016.
- [39] Alistair Knott. Moving towards responsible government use of ai in new zealand). <https://digitaltechip.nz/2021/03/22/moving-towards-responsible-government-use-of-ai-in-new-zealand/>.

- [40] I Elizabeth Kumar, Suresh Venkatasubramanian, Carlos Scheidegger, and Sorelle Friedler. Problems with shapley-value-based explanations as feature importance measures. In *International Conference on Machine Learning*, pages 5491–5500. PMLR, 2020.
- [41] Matt J Kusner, Joshua R Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. *arXiv preprint arXiv:1703.06856*, 2017.
- [42] Lydia T Liu, Sarah Dean, Esther Rolf, Max Simchowitz, and Moritz Hardt. Delayed impact of fair machine learning. In *International Conference on Machine Learning*, pages 3150–3158. PMLR, 2018.
- [43] Joshua R Loftus, Chris Russell, Matt J Kusner, and Ricardo Silva. Causal reasoning for algorithmic fairness. *arXiv preprint arXiv:1805.05859*, 2018.
- [44] Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*, 2017.
- [45] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>.
- [46] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635*, 2019.
- [47] Aditya Krishna Menon and Robert C Williamson. The cost of fairness in binary classification. In *Conference on Fairness, Accountability and Transparency*, pages 107–118, 2018.
- [48] Alexey Miroshnikov, Konstantinos Kotsiopoulos, Ryan Franks, and Arjun Ravi Kannan. Wasserstein-based fairness interpretability framework for machine learning models. *arXiv preprint arXiv:2011.03156*, 2020.
- [49] Razieh Nabi and Ilya Shpitser. Fair inference on outcomes. In *Proceedings of the... AAAI Conference on Artificial Intelligence. AAAI Conference on Artificial Intelligence*, volume 2018, page 1931. NIH Public Access, 2018.
- [50] Vedant Nanda, Samuel Dooley, Sahil Singla, Soheil Feizi, and John P Dickerson. Fairness through robustness: Investigating robustness disparity in deep learning. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 466–477, 2021.
- [51] Preetam Nandy, Cyrus Diccio, Divya Venugopalan, Heloise Logan, Kinjal Basu, and Noureddine El Karoui. Achieving fairness via post-processing in web-scale recommender systems, 2021.
- [52] Aravind Narayanan. 21 fairness definitions and their politics. <https://fairmlbook.org/tutorial2.html>.
- [53] K. Nishida, S. Shimada, S. Ishikawa, and K. Yamauchi. Detecting sudden concept drift with knowledge of human behavior. In *2008 IEEE International Conference on Systems, Man and Cybernetics*, pages 3261–3267, 2008. doi: 10.1109/ICSMC.2008.4811799.
- [54] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.
- [55] Government of Canada. Responsible use of artificial intelligence (ai). <https://www.canada.ca/en/government/system/digital-government/digital-government-innovations/responsible-use-ai.html>.
- [56] Osonde A Osoba and William Welser IV. *An intelligence in our image: The risks of bias and errors in artificial intelligence*. Rand Corporation, 2017.
- [57] Fábio Pinto, Marco OP Sampaio, and Pedro Bizarro. Automatic model monitoring for data streams. *arXiv preprint arXiv:1908.04240*, 2019.
- [58] Edward Raff, Jared Sylvester, and Steven Mills. Fair forests: Regularized tree induction to minimize model bias. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 243–250, 2018.
- [59] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

- [60] Marcos Salganicoff. Tolerating concept and sampling shift in lazy learning using prediction error context switching. In *Lazy learning*, pages 133–155. Springer, 1997.
- [61] Nithya Sambasivan, Erin Arnesen, Ben Hutchinson, Tulsee Doshi, and Vinodkumar Prabhakaran. Re-imagining algorithmic fairness in india and beyond. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 315–328, 2021.
- [62] Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet Vertesi. Fairness and abstraction in sociotechnical systems. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 59–68, 2019.
- [63] Aalok Shanbhag, Avijit Ghosh, and Josh Rubin. Unified shapley framework to explain prediction drift. *arXiv preprint arXiv:2102.07862*, 2021.
- [64] Ashudeep Singh and Thorsten Joachims. Fairness of exposure in rankings. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2219–2228, 2018.
- [65] Kenneth O Stanley. Learning concept drift with a committee of decision trees. *Informe técnico: UT-AI-TR-03-302, Department of Computer Sciences, University of Texas at Austin, USA*, 2003.
- [66] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pages 3319–3328. PMLR, 2017.
- [67] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Bias preservation in machine learning: The legality of fairness metrics under eu non-discrimination law. *West Virginia Law Review, Forthcoming*, 2021.
- [68] James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viégas, and Jimbo Wilson. The what-if tool: Interactive probing of machine learning models. *IEEE transactions on visualization and computer graphics*, 26(1):56–65, 2019.
- [69] Linda F Wightman. Lsac national longitudinal bar passage study. Isac research report series. 1998.
- [70] Steven L Willborn. The disparate impact model of discrimination: Theory and limits. *Am. UL Rev.*, 34: 799, 1984.
- [71] Christo Wilson, Avijit Ghosh, Shan Jiang, Alan Mislove, Lewis Baker, Janelle Szary, Kelly Trindel, and Frida Polli. Building and auditing fair algorithms: A case study in candidate screening. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 666–677, 2021.
- [72] Meike Zehlike, Francesco Bonchi, Carlos Castillo, Sara Hajian, Mohamed Megahed, and Ricardo Baeza-Yates. Fa\* ir: A fair top-k ranking algorithm. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 1569–1578, 2017.
- [73] Indre Žliobaite. Change with delayed labeling: When is it detectable? In *2010 IEEE International Conference on Data Mining Workshops*, pages 843–850. IEEE, 2010.
- [74] Indrė Žliobaitė. Combining similarity in time and space for training set formation under concept drift. *Intelligent Data Analysis*, 15(4):589–611, 2011.

## 10 Appendix

### 10.1 Proofs

Given that the attribution method A satisfies the efficiency axiom,  $QDD_b = \sum_{f=1}^F QDDA_{b,A,f}$ .

The QDD Attribution for feature f, for prediction sample  $S_1$  over  $S_2$  in bin b,  $QDDA_{b,A,f}$  is a measure of the change in QDD in bin b that can be attributed to feature f using attribution method A that satisfies the efficiency axiom.

$$QDDA_{b,A,f} = \frac{1}{N_t} \sum_{n=1}^{N_t} attr_{n,A,S_1,f} - \frac{1}{N_r} \sum_{n=1}^{N_r} attr_{n,A,S_2,f} \quad (4)$$

Since the attribution method A satisfies Efficiency, for each instance in the sample  $S_1$  and  $S_2$ ,  $\sum_{f=1}^F attributions_f = prediction - baseline prediction$

For the same baseline,

$$\begin{aligned} \therefore \frac{1}{N_t} \sum_{n=1}^{N_t} \sum_{f=1}^F attr_{n,A,S_1,f} - \frac{1}{N_r} \sum_{n=1}^{N_r} \sum_{f=1}^F attr_{n,A,S_2,f} &= \frac{1}{N_1} \sum_{n=1}^{N_1} S_{1,n} - \frac{1}{N_2} \sum_{n=1}^{N_2} S_{2,n} \\ \therefore QDD_b &= \sum_{f=1}^F QDDA_{b,A,f} \end{aligned}$$

### 10.2 Case Studies

1. The Luxembourg dataset [74] is constructed using European Social Survey data. Each instance is an individual. The labels indicate high or low internet usage, which has a continuous output. We trained a five layer neural network regression model using the following features:

- Country - cntry
- Total time reading the news on average weekday -nwspol
- Personal use of internet/e-mail/www - netusm
- How interested in politics - polintr
- Trust in the European Parliament - trstep
- Trust in the United Nations - trstun
- Signed petition last 12 months - sgnptit
- Member of political party - wrkppty
- How happy are you - happy
- How often socially meet with friends, relatives or colleagues - sclmeet
- Subjective general health - health
- How often pray apart from at religious services - pray
- How long ago first came to live in country - livecnta
- Number of people living regularly as member of household - hhmmb
- Second person in household:Relationship to respondent - rshipa2
- Highest level of education - edulvlb
- Main activity, last 7 days. All respondents. Post coded - mnactic
- Responsible for supervising other employees - jbspv

2. The LSAC dataset [69] is a public dataset containing information about students in US Law School admissions process. We trained a two layer neural network model using the following features:

- Law School SAT score - lsat
- Law School GPA - gpa
- Undergraduate GPA - ugpa

### **10.3 Data Ethics**

Both datasets used in the case studies are publicly available datasets and contain no personally identifying information. They do not therefore constitute Human Subject Research and are exempt from an IRB requirement according to US Code 45 C.F.R. §46.101.