

566 **A Disparity Metric Definitions**

567 **A.1 Observational Metrics**

568 **False Positive Rate Parity** Definition: $\hat{Y} \perp A \mid Y = 0$

569 Measured as: $P(\hat{Y} = 1 \mid A = 0, Y = 0) - P(\hat{Y} = 1 \mid A = 1, Y = 0)$

570 **False Negative Rate Parity** Definition: $\hat{Y} \perp A \mid Y = 1$

571 Measured as: $P(\hat{Y} = 1 \mid A = 0, Y = 1) - P(\hat{Y} = 1 \mid A = 1, Y = 1)$

572 **Positive Predictive Parity** Definition: $Y \perp A \mid \hat{Y} = 1$

573 Measured as: $P(Y = 1 \mid A = 0, \hat{Y} = 1) - P(Y = 1 \mid A = 1, \hat{Y} = 1)$

574 **Negative Predictive Parity** Definition: $Y \perp A \mid \hat{Y} = 0$

575 Measured as: $P(Y = 1 \mid A = 0, \hat{Y} = 1) - P(Y = 1 \mid A = 1, \hat{Y} = 0)$

576 **Equalized Odds** Definition: $Y \perp A \mid \hat{Y}$

577 Measured as: $\max \left\{ \text{FPR}(Y, A, \hat{Y}), \text{FNR}(Y, A, \hat{Y}) \right\}$ for false positive rate (FPR) and false negative
578 rate (FNR) given above.

579 **A.2 ECP Parity Metric Definitions**

580 **Counterfactual False Positive Rate Parity** Definition: $\hat{Y} \perp A \mid Y(D = 1) = 0$

581 Measured as: $P_C(\hat{Y} = 1 \mid A = 0, Y(D = 1) = 0) - P(\hat{Y} = 1 \mid A = 1, Y(D = 1) = 0)$

582 **Counterfactual False Negative Rate Parity** Definition: $\hat{Y} \perp A \mid Y = 1$

583 Measured as: $P(\hat{Y} = 1 \mid A = 0, Y(D = 1) = 1) - P(\hat{Y} = 1 \mid A = 1, Y(D = 1) = 1)$

584 **Counterfactual Positive Predictive Parity** Definition: $Y(D = 1) \perp A \mid \hat{Y} = 1$

585 Measured as: $P(Y(D = 1) = 1 \mid A = 0, \hat{Y} = 1) - P(Y(D = 1) = 1 \mid A = 1, \hat{Y} = 1)$

586 **Counterfactual Negative Predictive Parity** Definition: $Y(D = 1) \perp A \mid \hat{Y} = 0$

587 Measured as: $P(Y(D = 1) = 1 \mid A = 0, \hat{Y} = 1) - P(Y(D = 1) = 1 \mid A = 1, \hat{Y} = 0)$

588 **Counterfactual Equalised Odds** Definition: $Y(D = 1) \perp A \mid \hat{Y}$

589 Measured as: $\max \left\{ \text{CF}_F\text{PR}(Y, A, \hat{Y}), \text{CF}_F\text{NR}(Y, A, \hat{Y}) \right\}$ for counterfactual false positive rate
590 (CF_FPR) and counterfactual false negative rate (CF_FNR) given above.

591 **B Technical Description**

592 **B.1 Marginalisation in DAGs**

593 **Marginalisation Operation** Suppose \mathbf{V} can be split as $\mathbf{V} = \tilde{\mathbf{V}} \cup \tilde{\mathbf{U}}$ where we are interested in the
594 causal structure over $\tilde{\mathbf{V}}$ and do not observe the variables $\tilde{\mathbf{U}}$. We start from a causal graph \mathcal{G} , with
595 unobserved $\tilde{\mathbf{U}}$ we marginalise to get to a graph \mathcal{G}' which is of the form of Definition [1](#) by doing the
596 following:

- 597 1. For all $U \in \tilde{\mathbf{U}}$, add an edge $Z \rightarrow \tilde{Z}$ if the current graph contains $Z \rightarrow U \rightarrow \tilde{Z}$ and then
598 delete any edges $Z \rightarrow U$,
- 599 2. After completing the first step for all variables in $\tilde{\mathbf{U}}$, delete any U if there exists another
600 $\tilde{U} \in \tilde{\mathbf{U}}$ that influences all of the variables U influences.

601 Evans [\[22\]](#) showed that there is a structural causal model over the resulting graph which preserves the
602 causal structure over the variables $\tilde{\mathbf{V}}$. Importantly, due to the deletion step, this model has bounded
603 number of unobserved variables, regardless of how large the set $\tilde{\mathbf{U}}$ is.

604 **Graphical examples**

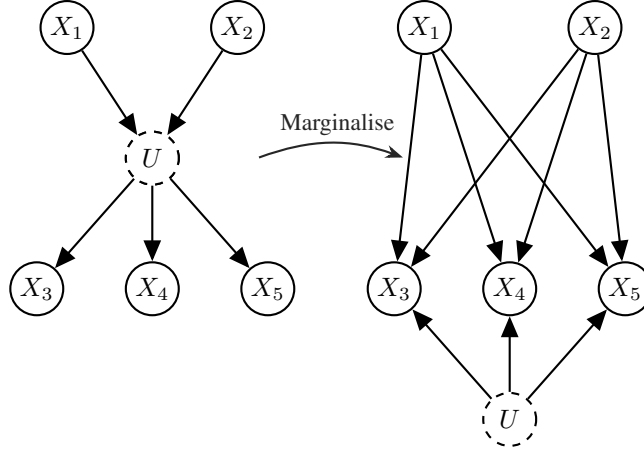


Figure 5: Example of step one in the marginalisation, taken from Evans [22].

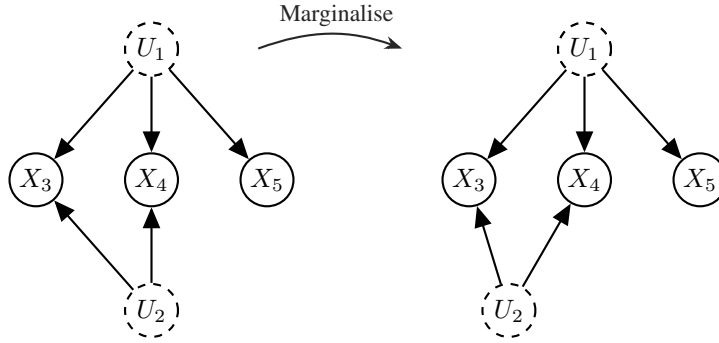


Figure 6: Example of step two in the marginalisation.

605 B.2 Alternative Causal Graphs for Proxy Bias

606 Here we provide the following result demonstrating that a wide variety of Graphs can give the same
 607 outcome under proxy bias:

608 **Proposition 1.** *So long as any additional unobserved variables U' satisfy the following:*

- 609 1. U' does not cause A .
- 610 2. There is no direct arrow from U' to \hat{Y} .

611 Then marginalising over U' will lead to the same graph as Figure 1a

612 *Proof.* To show this we need to demonstrate that once we have performed the marginalisation
 613 operations, no additional edges or nodes will be added to the graph. We do this step by step:

- 614 1. This step will add edges if we have two vertices V, V' such that $V \rightarrow U' \rightarrow V'$. However,
 615 if neither of V, V' are \hat{Y} then these vertices will already be adjacent in the graph. As the
 616 graph is acyclic that means we cannot be adding any edges.
- 617 2. After removing all edges in step one, we will be left so that U' has no parents and effects a
 618 subset of vertices in the graph. However, as U' does not cause A , this must be a subset of
 619 $\{\hat{Y}, Y_P, Y\}$. As these are the vertices caused by U this will lead to the deletion of U' .

620 □

Dataset	Task	Proxy Bias	Selection Bias	ECP Bias
Adult	Synthetic			
KDD Census-Income	Synthetic			
German credit	Credit risk		✓	✓
Dutch census	Synthetic			
Bank marketing	Client		✓	
Credit card clients	Default Risk		✓	✓
COMPAS recid.	Risk prediction	✓	✓	✓
COMPAS viol. recid.	Risk prediction	✓	✓	✓
Communities&Crime	Neighborhood risk	✓	✓	✓
Diabetes	Re-admission risk	✓		✓
Ricci	Promotion Prediction	✓	✓	✓
Student-Mathematics	Admissions	✓	✓	✓
Student-Portuguese	Admissions	✓	✓	✓
OULAD	Admissions	✓	✓	✓
Law School	Admissions	✓	✓	✓

Table 2: Analysis of the datasets from Le Quy et al. [43], split by task. The explanation for the biases are given in Appendix E.

621 C Cross Dataset Analysis

622 In this section we analyse the datasets presented in Le Quy et al. [43] for the three biases we present
623 in Section 3. We describe each dataset, give the task which most closely relates to the use of this
624 dataset, and relative to this task we decide if each of the three measurement biases are present or not.
625 For each bias we provide a justification of our decision.

626 **Synthetic tasks** The synthetic tasks are hard to discuss since the biases are contextual and these
627 tasks are purely theoretical. Given a downstream task they might or might not have the biases we
628 discuss. Therefore we drop them from the analysis.

629 **Bank marketing Dataset** The goal here is to target current clients for the bank to open more
630 accounts. Since the outcome in this case is exactly what the bank seeks to maximise, this dataset
631 does not exhibit proxy or ECP bias. However, contacts we made via phone, so there is selection bias
632 in whether people answered the phone.

633 **German credit and Credit card clients** For both of these datasets goal is to predict whether
634 customers face default risk. The aim is to use this to decide if applying customers present a risk to the
635 bank or not. As a result of this, there will be selection bias due to the fact that since defaults are only
636 observed for the firms’ previous customers. Finally as with the example in the main text, this exhibits
637 extra-classificatory policy bias since the firm sets the credit limit which impacts the likelihood of
638 default.

639 **COMPAS recid. and COMPAS viol. recid. and Communities and Crime** Datasets build off
640 COMPAS have been well documented to exhibit all these biases and more [5]. These issues are
641 not unique to COMPAS and are exhibited in all other recidivism and crime prediction datasets, as
642 such they will also apply to communities and crime, where the aim is to predict number of historical
643 crimes per hundred thousand population for a number of states. Moreover, a large degree of missing
644 values in this dataset show the issues due to selection bias.

645 **Diabetes** For this dataset, the goal is to predict if a patient will be readmitted in the next 30 days.
646 The aim is to use this to decide how much of a health risk a given patient is upon leaving the hospital,
647 to decide if they should be kept there. The population is a sample of the patient pool, and so there
648 should not be selection bias. Readmissions are different from the underling recurring illness, so this
649 does represent a proxy, albeit it a fairly reasonable one. In this case ECP bias is a cause for concern
650 due to the differences in quality of care by demographic group [21].

651 **Ricci** The Ricci dataset is an employment dataset, where the goal is to predict the likelihood of a
652 promotion based off of a selection of available covariates. A model trained on this data would then be

653 used to predict the potential of applying candidates in order to decide if they are invited to interview
 654 or do additional tests. This application would fall risk of all the biases we have presented and as such
 655 strong justification would be required as to the usefulness of the model. Going through one by one,
 656 proxy bias is exhibited in a similar way to the example presented in the main text, selection bias is
 657 present as the model is evaluated on a different population to the one it is trained on, and finally the
 658 firms policies will have an impact on who succeeds and is promoted at the company.

659 **Admissions datasets** The final datasets can all be grouped under admissions to academic institu-
 660 tions. Similarly to the employment example, these will exhibit all the biases we have outlined. This
 661 is because of the challenges of having a perfectly objective measure of performance, models being
 662 used on applying populations but fit on accepted populations, and the universities policies affecting
 663 the success of students. Therefore, when using these predictors arguments should be made about why
 664 using such a measure would not induce demographic skew.

665 D Additional Results

666 D.1 Proxy Label Results

667 D.1.1 Plots from Fogliato et al. [26] under varying assumptions

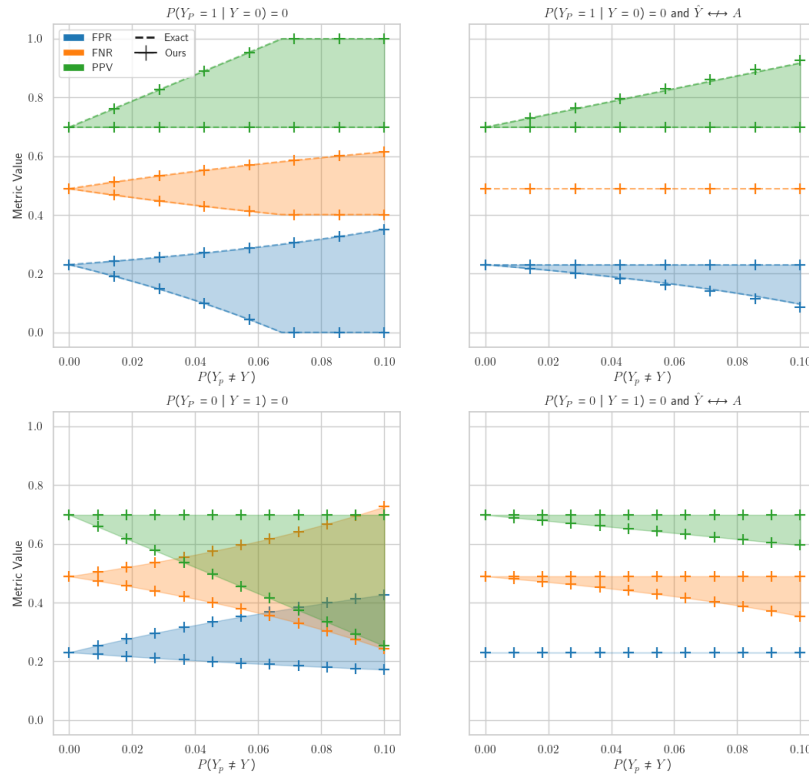


Figure 7: In this plot, we recreate the results from Fogliato et al. [26], where we are interested in the false positive rate (FPR), false negative rate (FNR), and positive predictive value (PPV) for a classifier trained on the COMPAS dataset. In this plot we consider varying for which j we have $P(Y_P = 1 - j | Y = j)$, and we can see that doing so greatly changes the shape of the sensitivity set. Moreover, when we pair these assumptions dropping of the red dashed edge in Figure 1a we see we can identify some of the metrics of interest under any degree of bias. For $j = 1$ we identify the FNR and for $j = 0$ we identify the FPR. We prove these identification results in Appendix D.1.2

668 **D.1.2 Proxy Identification Results**

669 In the set up in Fogliato et al. [26], the aim is the false positive/negative rate in a group $A = a$, where
 670 it is assumed that $P(Y = 1, Y_P = 0) = 0$. Now declaring the following parameters:

$$\begin{aligned} p_{ij} &= P(Y_P = i, \hat{Y} = j \mid A = a) \\ \alpha_j &= P(Y = 1, Y_P = 0, \hat{Y} = j \mid A = a) \\ \alpha &= \alpha_0 + \alpha_1 \end{aligned}$$

671 Under these assumptions α_0, α_1 are sufficient to parameterise the distribution, $P(Y, Y_P, \hat{Y} \mid A = a)$.
 672 Now, following [26] we have that:

$$\begin{aligned} \text{FPR}_Y &= \frac{p_{01} - \alpha_1}{p_{00} + p_{01} - \alpha} \\ \text{FNR}_Y &= \frac{p_{10} + \alpha_0}{p_{10} + p_{11} + \alpha} \\ \text{PPV}_Y &= \frac{p_{11} + \alpha_1}{p_{01} + p_{11}} \end{aligned}$$

673 Now, with the absence of the dashed edge, the DAG in Figure 1a implies the independence $\hat{Y} \perp Y_P \mid$
 674 Y, A . Therefore we get the following:

$$\begin{aligned} \alpha_j &= P(Y = 1, Y_P = 0, \hat{Y} = j \mid A = a) \\ &= \frac{P(Y = 1, Y_P = 0 \mid A = a)P(Y = 1, \hat{Y} = j \mid A = a)}{P(Y = 1 \mid A = a)} \\ &= \frac{\alpha(p_{1j} + \alpha_j)}{p_{10} + p_{11} + \alpha} \end{aligned}$$

675 Solving for α_j , we get $\alpha_j = \alpha \left(\frac{p_{1j}}{p_{10} + p_{11}} \right)$. Now, inputting this for α_0 in the expression for FNR_Y
 676 we get:

$$\begin{aligned} \text{FNR}_Y &= p_{10} \left(\frac{1 + \frac{\alpha}{p_{10} + p_{11}}}{p_{10} + p_{11} + \alpha} \right) \\ &= \frac{p_{10}}{p_{10} + p_{11}} \\ &= \text{FNR}_{Y_P} \end{aligned}$$

677 Therefore, under the assumptions given, the true false negative rate is identified and equal to the
 678 observed false negative rate on the proxy labels. Inputting the value for α_1 into FPR_Y we instead
 679 get:

$$\text{FPR}_Y = \frac{p_{01} - \alpha \left(\frac{p_{10}}{p_{10} + p_{11}} \right)}{(p_{00} + p_{01} - \alpha)}$$

680 As this is a decreasing function of α we can see that for $\alpha \leq \alpha_0$, FPR_Y is bounded as:

$$\frac{p_{01}}{(p_{00} + p_{01})} \leq \text{FPR}_Y \leq \frac{p_{01} - \alpha \left(\frac{p_{10}}{p_{10} + p_{11}} \right)}{(p_{00} + p_{01} - \alpha)}$$

681 For PPV, we again input α_1 to give:

$$\text{PPV}_Y = \frac{p_{11} + \alpha \left(\frac{p_{11}}{p_{10} + p_{11}} \right)}{p_{01} + p_{11}}$$

682 Leading to the bounds:

$$\text{PPV}_{Y_P} \leq \text{PPV}_Y \leq \frac{p_{11} + \alpha \left(\frac{p_{11}}{p_{10} + p_{11}} \right)}{p_{01} + p_{11}}$$

683 The statements for the identification of the false positive rate and false negative rate are as follows:

684 **Proposition 2.** Suppose we have $P(Y_P = 1 | Y = 0) = 0$. Then under the conditional independence
 685 statement $\hat{Y} \perp Y_P | Y, A$, for all level of proxy bias $P(Y_P \neq Y)$:

$$\text{FNR}_{Y|A=a} = \text{FNR}_{Y_P|A=a}$$

686 Where $\text{FNR}_{Y|A=a}$ is the true false negative rate for the group $A = a$ and $\text{FNR}_{Y_P|A=a}$ is the proxied
 687 false negative rate.

688 *Proof.* Follows from the above derivations. □

689 Now the equivalent statement for the false positive ratio:

690 **Proposition 3.** Suppose we have $P(Y_P = 0 | Y = 1) = 0$. Then under the conditional independence
 691 statement $\hat{Y} \perp Y_P | Y, A$, for all level of proxy bias $P(Y_P \neq Y)$:

$$\text{FPR}_{Y|A=a} = \text{FPR}_{Y_P|A=a}$$

692 Where $\text{FPR}_{Y|A=a}$ is the true false negative rate for the group $A = a$ and $\text{FPR}_{Y_P|A=a}$ is the proxied
 693 false negative rate.

694 *Proof.* This follows from considering the distribution where Y, Y_P and \hat{Y} are all flipped as any
 695 statement about the false positive rate in the original distribution translates to a statement about
 696 the false negative rate in the flipped distribution. The assumption $P(Y_P = 0 | Y = 1)$ in the
 697 original distribution translates to $P(Y_P = 1 | Y = 0)$ in the flipped distribution, whereas all other
 698 assumptions are symmetric to the flipping operation. Therefore we can apply proposition 2 to see
 699 that the flipped FNR is constant under any degree of proxy noise. This leads us to conclude that
 700 under these assumptions the FPR in the original distribution must also be constant under any degree
 701 of proxy noise. □

702 D.2 Selection Results

703 D.2.1 Selective labels under MNAR

704 Here we include an experiment applying the framework to selective labels under the missing not a
 705 random assumption (MNAR) [56]. This supposes that we only see the outcome on a subset of the full
 706 dataset, with the outcome on the rest of the dataset free to vary arbitrarily. We work with the Dutch
 707 census dataset Van der Laan [60], first fitting an unconstrained logistic regression, then forming the
 708 selected population as those who have a predicted probability higher 0.3.

709 Once we have formed the selected subset we then train four classifiers, each to satisfy a different
 710 parity metric. We train to false negative rate parity, false positive rate parity, positive predictive parity
 711 and negative predictive parity. False negative/positive rate parity are trained using the reductions
 712 approach [2], whereas for positive predictive parity and negative predictive parity we train 100
 713 predictors, each weighting different parts of the distribution, taking the one with the lowest parity
 714 score above a given accuracy threshold. The plots are shown in Figure 8.

715 D.2.2 Selection and Proxy Plots

716 In this we demonstrate the effect of selection and proxy bias jointly on the adult dataset. We include
 717 the results in Figure 9, which show that both the occurrence of multiple biases acts differently for
 718 different parity metrics.

719 D.3 ECP bias results

720 D.3.1 ECP experimental set up

721 For this experiment, we focus on finding the possible ranges for the counterfactual parity metrics
 722 from the given observational statistics. We use the sensitivity parameter of $P(Y(1) \neq Y(0))$, adding
 723 additional causal assumptions such as monotonicity ($Y(1) \geq Y(0)$) and if the policy is observed or
 724 not. When simulating the policy we draw $\text{ECP} \sim \text{Ber}(\frac{1}{2} + c * A)$ for $c = 0.2$ in order to skew the
 725 policy in one direction. Results show in Figure 10.

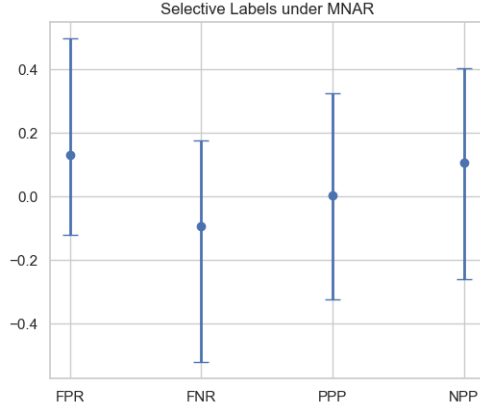


Figure 8: This plot demonstrates a sensitivity analysis for selective labels on the Dutch dataset under the missing not at random assumption.

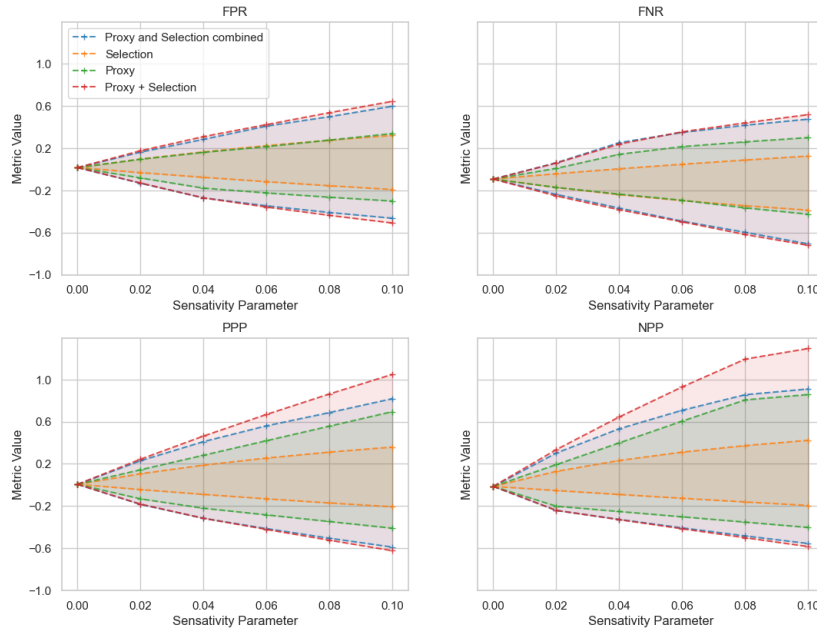


Figure 9: In these plots, we can see the effect of doing a sensitivity analysis jointly for selection and proxy bias. We can see that for the false positive rate parity (FPR) and false negative rate parity (FNR) the combined bias behaves roughly as the sum of both biases, however for positive predictive parity (PPP) and negative predictive parity (NPP) the combination behaves differently with the combined bias amounting a smaller possible range for the metrics than the sum of the range of both biases individually.

726 **D.4 Causal Fairness Experiments**

727 In this section we give some results on applying our sensitivity analysis framework to causal fairness
 728 metrics of the variety detailed in [49]. Before doing so, we add some technical comments on these
 729 types of interventions in FairML, and some nuances of measurement bias in the context of causal
 730 inference.

731 Firstly, we would like to comment that our framework can still be applied to perform sensitivity
 732 analysis for measurement bias for other fairness metrics *without* having to consider counterfactuals
 733 relative protected characteristics such as race or gender, thereby avoiding difficulties with intervention

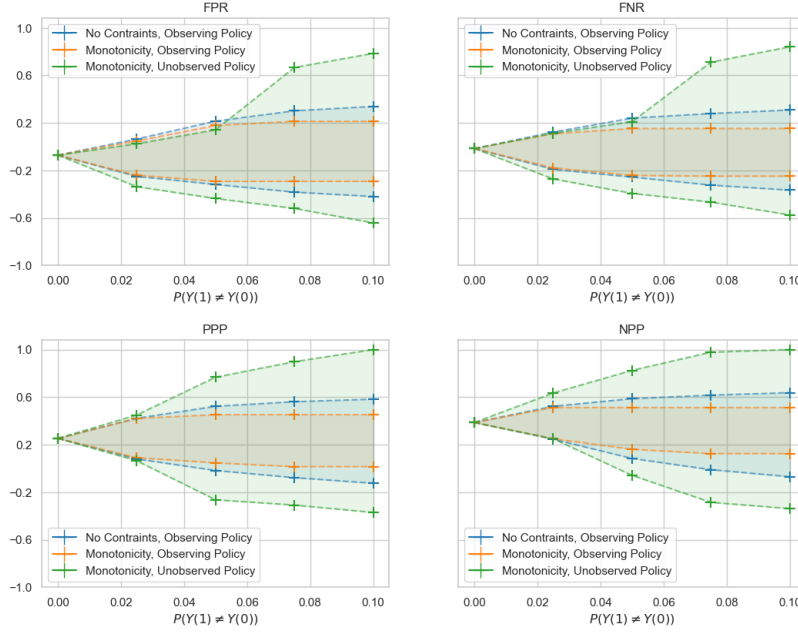


Figure 10: In these plots, we perform a sensitivity analysis for the value of the counterfactual parity metrics given in Appendix A.2. In each case we work under 3 differing levels of assumption and information

734 on such traits [39, 32, 36]. In this case A could be seen as denoting membership to a group and
 735 indexing different graphs for each group as in Bright et al. [10]. In this case the arrows leading from
 736 A would only express conditional independence relationships as opposed to causal ones. Notably in
 737 the graphs we suggest they are unconstrained.

738 Secondly, measurement biases and specifically selection bias in causal fairness comes with additional
 739 problems. This is because almost always, membership of such a dataset is causally downstream
 740 of the protected attribute, meaning that when conditioning on individuals being in a dataset we are
 741 introducing selection bias in some form. As Fawkes et al. [25] argue, this means DAG models will be
 742 unable to correctly capture the causal structure in most datasets we come across in FairML. Failing to
 743 account for such effects can lead to erroneous causal conclusions.

744 Having said this, we will proceed with applying the causal graphs in Figure 1 to do causal fairness
 745 analysis for the following metrics:

746 **Counterfactual Fairness (CF)** [41] We measure this as $P_C(\hat{Y}(A = 1) \neq \hat{Y}(A = 0))$ which is
 747 equal to 0 exactly when \hat{Y} is counterfactually fair [24].

748 **Total Effect (TE)** [49] Measured as $P_C(\hat{Y}(A = 1)) - P_C(\hat{Y}(A = 0))$.

749 **Spurious Effect (SE)** [49] Measured as $P_C(\hat{Y}(A = a)) - P_C(\hat{Y} | A = a)$.

750 Results are show in Figure 11, where we have assumed that counterfactual fairness is identified at
 751 a particular value. We can see that all causal fairness metrics recover a linear relationship under
 752 selection in this context.

753 E Details of cross dataset bias analysis

754 In this section we analyse the datasets presented in Le Quy et al. [43] for the three biases we present
 755 in Section 3. We describe each dataset, give the task which most closely relates to the use of this
 756 dataset, and relative to this task we decide if each of the three measurement biases are present or not.
 757 For each bias we provide a justification of our decision.

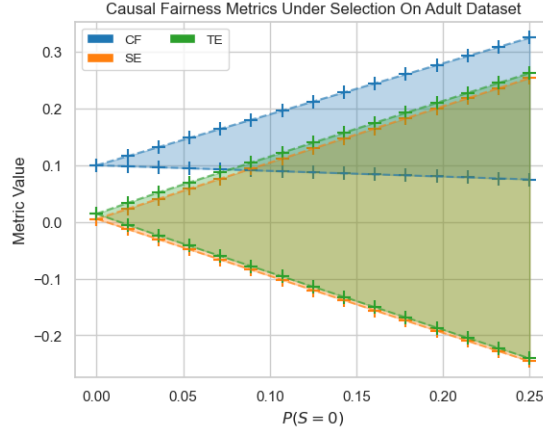


Figure 11: Causal fairness metrics under selection. We show plots for Counterfactual Fairness (CF), Total Effect (TE), and Spurious Effect (SE) using graph [1b](#) where we have additionally assumed that counterfactual fairness is point identified.

Dataset	Task	Proxy Bias	Selection Bias	ECP Bias
Adult	Synthetic			
KDD Census-Income	Synthetic			
German credit	Credit risk		✓	✓
Dutch census	Synthetic			
Bank marketing	Client		✓	
Credit card clients	Default Risk		✓	✓
COMPAS recid.	Risk prediction	✓	✓	✓
COMPAS viol. recid.	Risk prediction	✓	✓	✓
Communities&Crime	Neighborhood risk	✓	✓	✓
Diabetes	Re-admission risk	✓	✓	✓
Ricci	Promotion Prediction	✓	✓	✓
Student-Mathematics	Admissions	✓	✓	✓
Student-Portuguese	Admissions	✓	✓	✓
OULAD	Admissions	✓	✓	✓
Law School	Admissions	✓	✓	✓

Table 3: Analysis of the datasets from Le Quy et al. [\[43\]](#), split by task. The explanation for the biases are given in Appendix [E](#).

758 **Synthetic tasks** The synthetic tasks are hard to discuss since the biases are contextual and these
 759 tasks are purely theoretical. Given a downstream task they might or might not have the biases we
 760 discuss. Therefore we drop them from the analysis.

761 **Bank marketing Dataset** The goal here is to target current clients for the bank to open more
 762 accounts. Since the outcome in this case is exactly what the bank seeks to maximise, this dataset
 763 does not exhibit proxy or ECP bias. However, contacts we made via phone, so there is selection bias
 764 in whether people answered the phone.

765 **German credit and Credit card clients** For both of these datasets goal is to predict whether
 766 customers face default risk. The aim is to use this to decide if applying customers present a risk to the
 767 bank or not. As a result of this, there will be selection bias due to the fact that since defaults are only
 768 observed for the firms' previous customers. Finally as with the example in the main text, this exhibits
 769 extra-classificatory policy bias since the firm sets the credit limit which impacts the likelihood of
 770 default.

771 **COMPAS recid. and COMPAS viol. recid. and Communities and Crime** Datasets build off
 772 COMPAS have been well documented to exhibit all these biases and more [\[5\]](#). These issues are
 773 not unique to COMPAS and are exhibited in all other recidivism and crime prediction datasets, as

774 such they will also apply to communities and crime, where the aim is to predict number of historical
775 crimes per hundred thousand population for a number of states. Moreover, a large degree of missing
776 values in this dataset show the issues due to selection bias.

777 **Diabetes** For this dataset, the goal is to predict if a patient will be readmitted in the next 30 days.
778 The aim is to use this to decide how much of a health risk a given patient is upon leaving the hospital,
779 to decide if they should be kept there. The population is a sample of the patient pool, and so there
780 should not be selection bias. Readmissions are different from the underlying recurring illness, so this
781 does represent a proxy, albeit it a fairly reasonable one. In this case ECP bias is a cause for concern
782 due to the differences in quality of care by demographic group [21].

783 **Ricci** The Ricci dataset is an employment dataset, where the goal is to predict the likelihood of a
784 promotion based off of a selection of available covariates. A model trained on this data would then be
785 used to predict the potential of applying candidates in order to decide if they are invited to interview
786 or do additional tests. This application would fall risk of all the biases we have presented and as such
787 strong justification would be required as to the usefulness of the model. Going through one by one,
788 proxy bias is exhibited in a similar way to the example presented in the main text, selection bias is
789 present as the model is evaluated on a different population to the one it is trained on, and finally the
790 firms policies will have an impact on who succeeds and is promoted at the company.

791 **Admissions datasets** The final datasets can all be grouped under admissions to academic institu-
792 tions. Similarly to the employment example, these will exhibit all the biases we have outlined. This
793 is because of the challenges of having a perfectly objective measure of performance, models being
794 used on applying populations but fit on accepted populations, and the universities policies affecting
795 the success of students. Therefore, when using these predictors arguments should be made about why
796 using such a measure would not induce demographic skew.

797 **F Details of cross dataset experiment**

798 For this experiment, we train numerous predictors across a variety of common fairness benchmarking
799 datasets [43] to satisfy parity constraints. For each dataset we train 18 classifiers total, where the
800 model ML is one of logistic regression, naïve Bayes and a decision tree and the parity constraint is
801 false negative rate parity, false positive rate parity, positive predictive parity and negative predictive
802 parity, demographic parity and equalized odds. With the exception of positive/negative predictive
803 parity we train all classifier to satisfy these constraints using the reductions approach [2]. For
804 positive/negative predictive parity, we train 100 predictors, each weighting different parts of the
805 distribution, taking the one with the lowest parity score above a given accuracy threshold. We vary
806 the sensitivity parameter over a range of realistic values for many real-world settings, computing the
807 sensitivity bounds for each level of the parameter. We find that, except for demographic parity, all
808 parity measures we evaluate exhibit significant sensitivity over these parameter ranges. This makes
809 it hard to understand what satisfying, e.x. equalised odds means on a given dataset. The caveat is
810 that equalised odds is only satisfied as long as there are no significant measurement biases in the
811 underlying data, which is almost never the case in FairML audits.

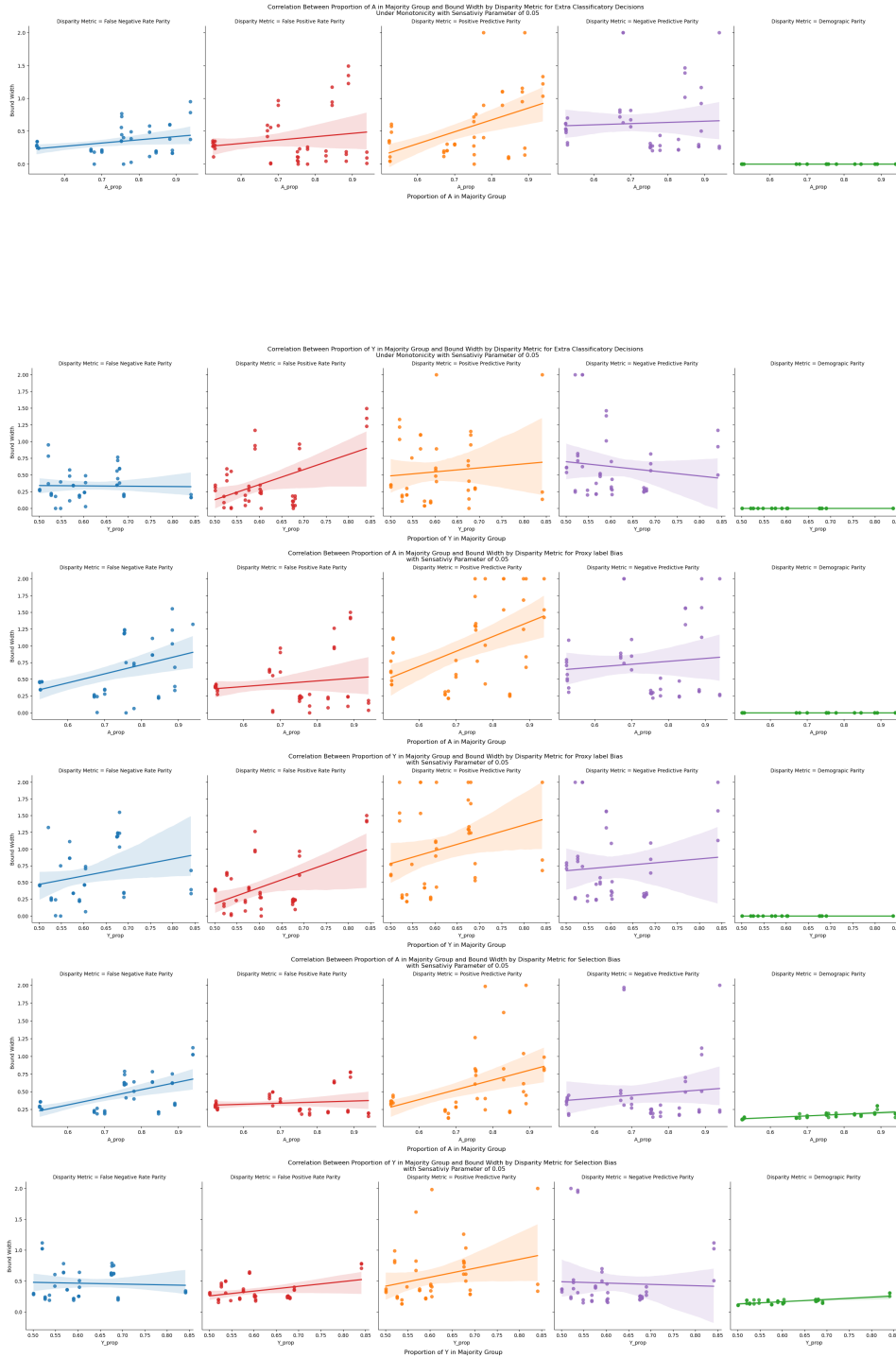
812 **F.1 Analysis of Results**

813 **F.1.1 Correlational Plots**

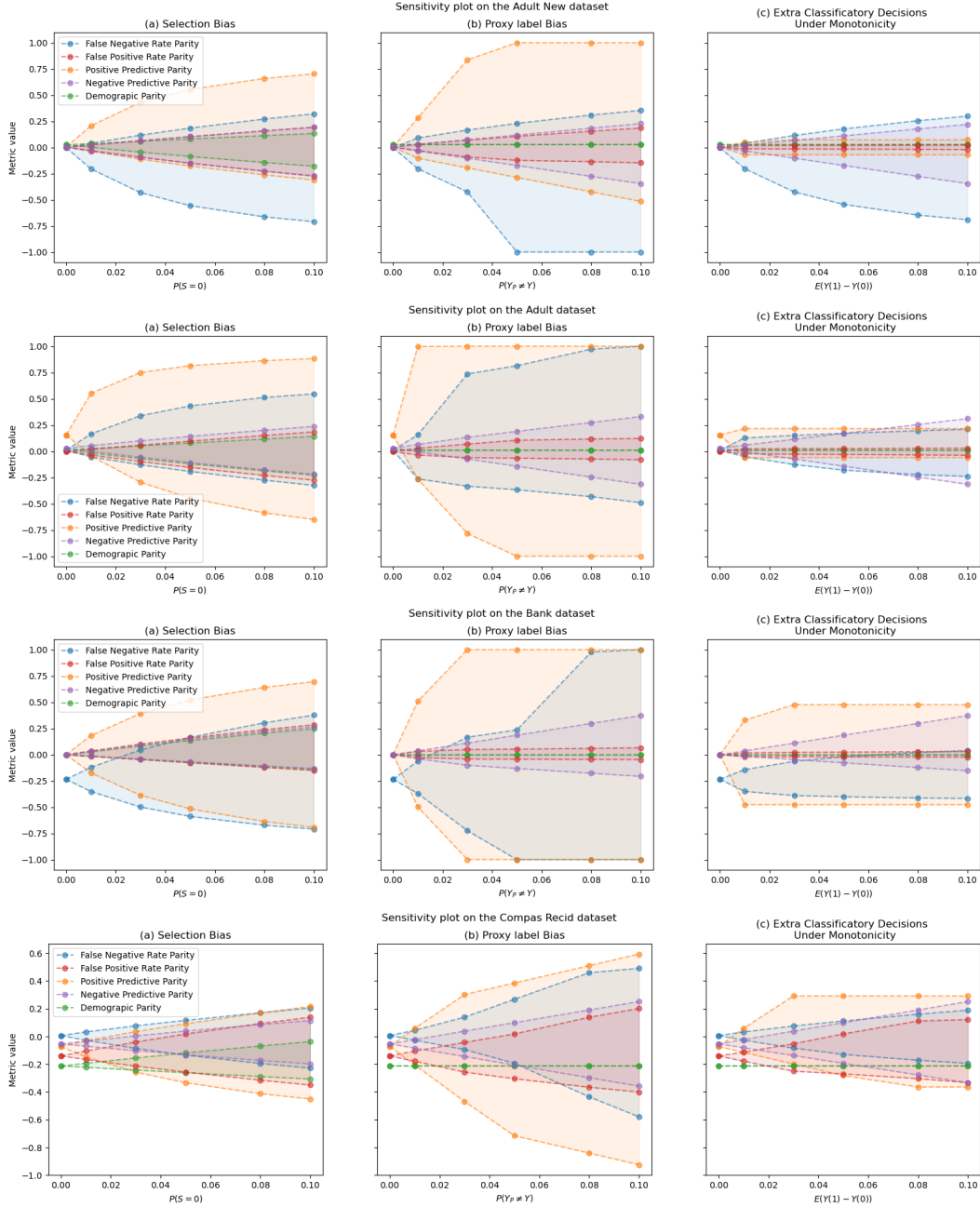
814 Here we explore how sensitivity varies according to class imbalance in either A or Y . We find that
815 for some metrics (Negative Predictive Parity) class imbalance seems to make little difference to the
816 sensitivity of metrics, with next to no correlation observed between imbalance and sensitivity. This
817 lies in contrast to other metrics (Positive Predictive Parity) where we can see a much more clear,
818 positive, correlation between class imbalance and sensitivity.

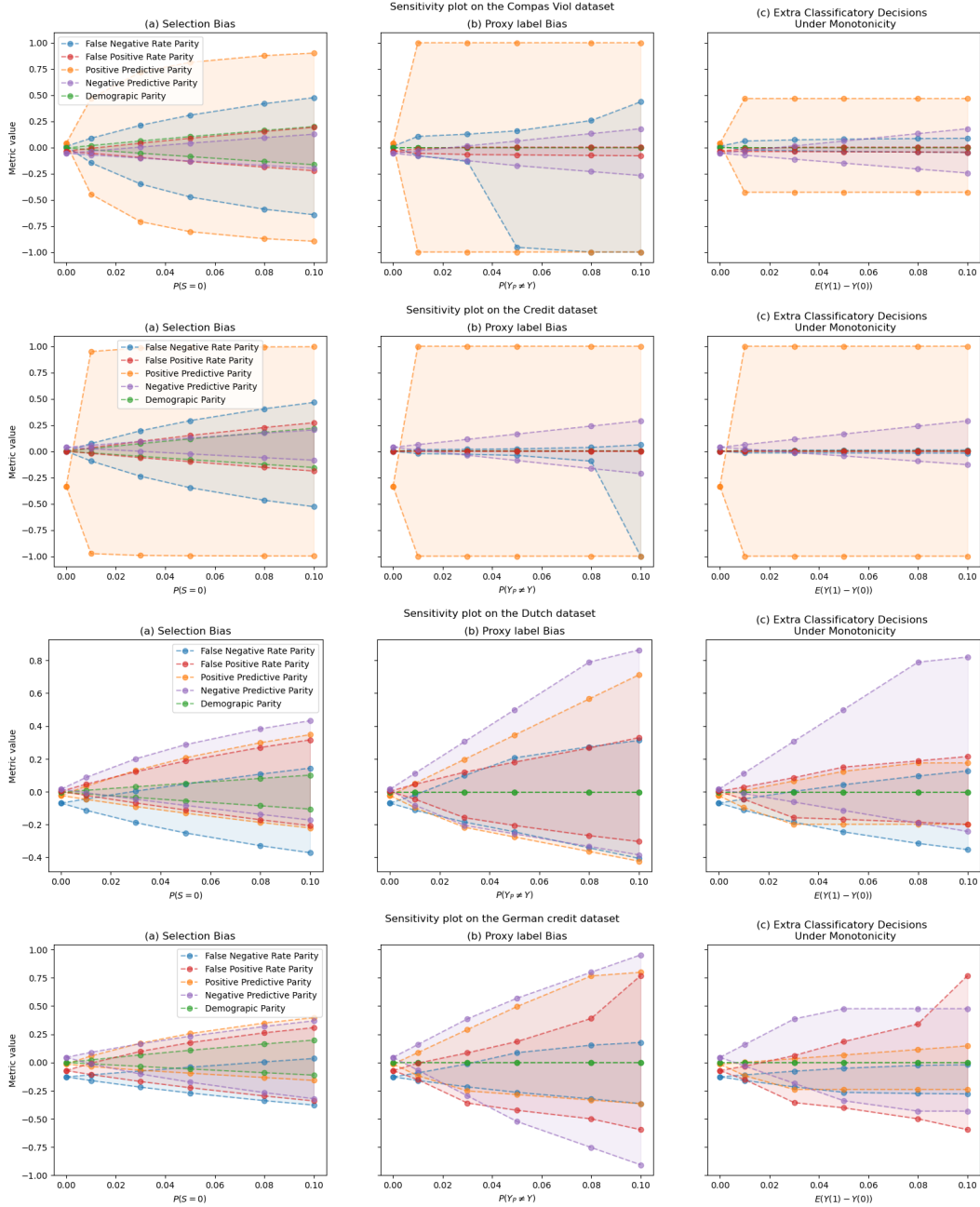
819 **F.1.2 Cross Dataset Analysis**

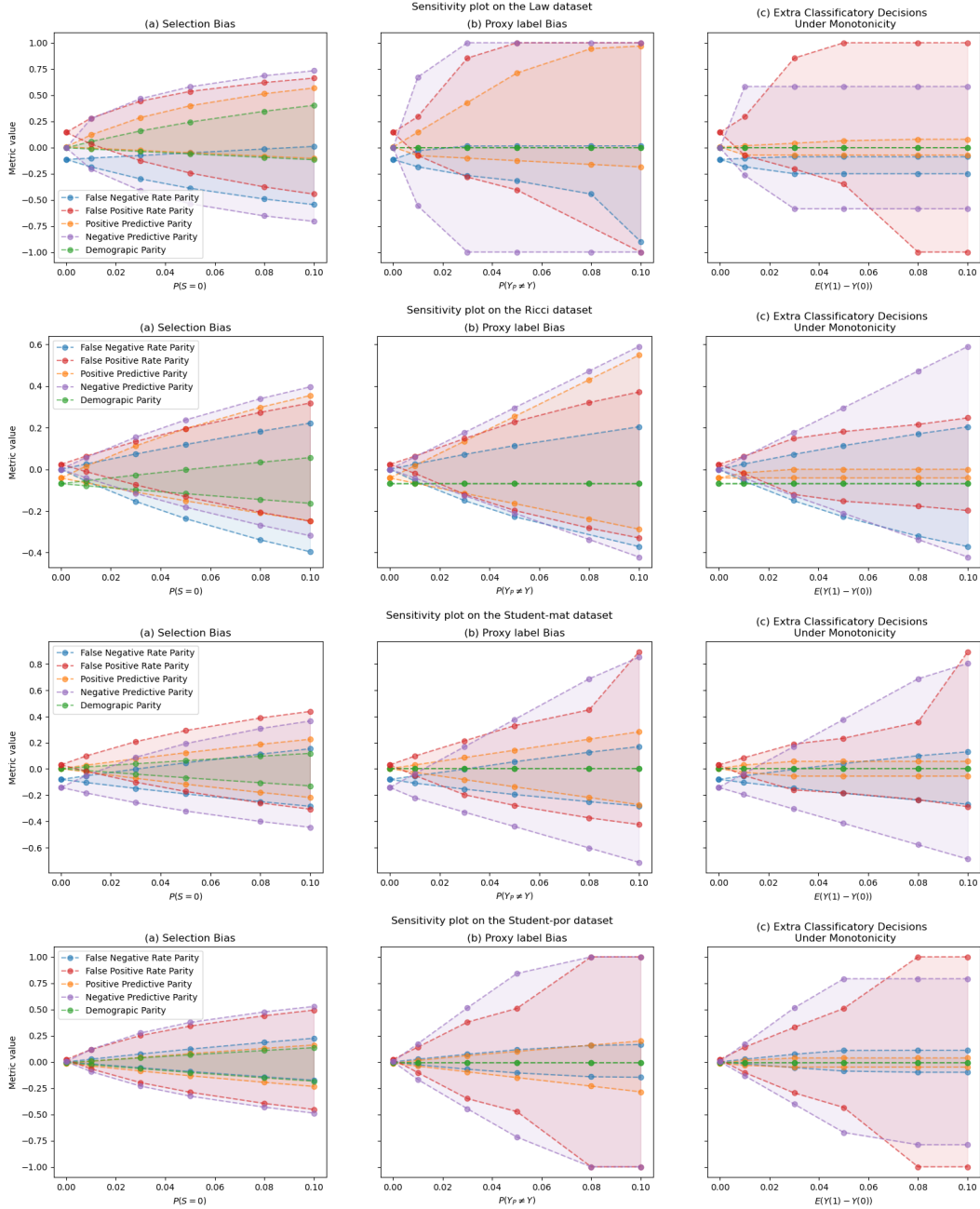
820 These results unpack some heterogeneity across datasets. We find that there are not hard and fast
821 rules here, each dataset and each bias requires its own analysis. At the same time this is broadly
822 consistent with our central contention that complexity goes hand in hand with fragility. In particular,



823 on the Adult New, Adult, Bank, Compas Recid/Viol, Credit, Dutch we see positive predictive parity is
 824 a standout fragile method across biases, followed by false negative rate pair and negative predictive
 825 parity. On the German Credit, Law, and Ricci, Student mat/por we see NPP and FNRP tend to be the
 826 worst, followed by FPRP and PPP.







827 **F.2 Impact Statement**

828 This work aims to broaden the discussion of measurement biases in FairML and provide practical
 829 tools for practitioners in the area to use. Our hope is that any potential societal consequences of the
 830 work will be positive, corresponding to more equitable algorithmic decision making.