





Beyond Predictive Accuracy: Fairness and Bias in Predicting Test Anxiety

Oscar Blessed Deho¹ , Srecko Joksimovic¹ , Maria Vieira¹ , and Ryan Baker² 

¹ Centre for Change and Complexity in Learning (C3L), University of South Australia, Adelaide SA 5001, Australia

`{oscar.deho, srecko.joksimovic, maria.vieira}@unisa.edu.au`

² Penn Center for Learning Analytics, University of Pennsylvania, Philadelphia, PA 19104, USA

`rybaker@upenn.edu`

Abstract. Test anxiety significantly impacts students’ academic performance and mental health, with complex interactions influenced by behavioral and demographic factors. This study examines the relationship between metacognitive self-regulation (MSR) behaviors and test anxiety across demographic groups, explores trade-off between predictive accuracy and fairness in test anxiety prediction models, and investigates how intersecting demographic attributes shape biases. The findings show that specific MSR behaviors, such as classroom distraction and frequent adaptation of study methods, are strongly correlated with test anxiety, highlighting key areas for targeted interventions. Demographic disparities are evident, with females experiencing higher levels of test anxiety and White students reporting more classroom distractions. A trade-off between predictive accuracy and fairness is observed, with highly accurate models not always performing well in terms of fairness, emphasizing the need for balanced model selection. Additionally, the study challenges traditional additive assumptions about fairness, finding that the intersection of demographic attributes produces unexpected compounded effects, such as compounded advantages for Non-White Migrants and mixed outcomes for White Females. We offer insights for designing accurate and equitable predictive models for test anxiety.

Keywords: Intersectional Fairness · Algorithmic Bias · Test Anxiety.

1 Introduction

Test anxiety is a widespread challenge for students, with well-documented effects on academic performance and mental health [4,45,39]. It disrupts cognitive functioning, reduces focus, and contributes to cycles of stress and underachievement [7,8]. Among many factors, metacognitive self-regulation (MSR)—the ability to plan, monitor, and adapt learning strategies—has been strongly linked to test anxiety [8]. However, this relationship is not straightforward. Demographic factors such as race, sex, and migration status may influence how students experience test anxiety and engage in MSR behaviors. For instance, some studies

show that female students, despite employing MSR strategies more frequently, still report higher levels of test anxiety than their male counterparts [39,8]. Understanding these variations is crucial for addressing disparities and designing interventions that are equitable and effective [30].

In recent years, predictive modeling plays a key role in identifying students at risk of test anxiety by analyzing behavioral, cognitive, and demographic data for early intervention [12]. By leveraging these insights, educators and researchers aim to target support toward students most in need. However, the implementation of predictive models raises important questions about fairness. Predictive systems that perform differently across demographic groups may unintentionally exacerbate inequalities rather than reduce them [2,16,32,15]. For example, if a model underestimates the risk of test anxiety for certain groups—such as Non-White students or Female Migrants—it could result in inadequate support for these populations [2,32]. While fairness in predictive modeling has received attention in educational research, many studies focus on individual demographic attributes, such as race or sex, without considering the compounded effects of intersecting identities [44,32]. For instance, being both a female and a migrant may create unique vulnerabilities that are not captured by single-attribute fairness evaluations [11]. Furthermore, fairness is often treated as separate from model performance, leaving the trade-off between *predictive accuracy* and fairness largely under-explored [19,43,48]. **Note:** In this paper, *predictive accuracy* refers to a model’s performance evaluated using one or more metrics such as precision and accuracy. In contrast, the *accuracy* metric specifically denotes the ratio of correct predictions to the total number of predictions. Investigating the trade-off between predictive accuracy and fairness is crucial in test anxiety prediction, as it affects the system’s utility for vulnerable groups [24].

Despite the progress made in understanding test anxiety and the use of predictive modeling in education, significant gaps remain. The relationship between MSR behaviors and test anxiety, particularly across demographic groups, is not well understood. Additionally, fairness considerations in predictive models for test anxiety require more attention, especially in balancing the trade-off between predictive accuracy and fairness. Lastly, evaluations of fairness must move beyond isolated attributes to address how intersecting demographic factors shape outcomes. To address these challenges, this study explores the following research questions (RQs):

RQ1a: What is the relationship between MSR behaviors and test anxiety?

RQ1b: How do MSR behaviors and test anxiety differ across demographic groups such as race, sex, and migrant status?

RQ2: Are the most *predictively accurate* models for test anxiety prediction also the most fair, or do trade-off exist between predictive accuracy and fairness?

RQ3: How does the intersection of demographic attributes compound or mitigate biases in predictive models for test anxiety?

We address the RQs by first using Spearman’s rank correlation and Mann-Whitney U tests to examine links between demographics, MSR behaviors, and test anxiety. Next, we train and evaluate five predictive models using four accu-

racy metrics and assess fairness across sex, migrant status, and race. Finally, we analyze the impact of intersecting demographics on fairness using a novel diagnostic metric. Our **contributions** are threefold: (1) we introduce the Residual Fairness Gap (RFG), a metric for assessing intersectional fairness in predictive models; (2) we demonstrate the trade-off between predictive accuracy and fairness, emphasizing the importance of informed model selection; and (3) we show that the combined effects of intersecting demographic attributes are often complex, going beyond simple additive assumptions.

2 Related Works

2.1 MSR and Test Anxiety

The relationship between test anxiety and MSR is a topic of ongoing discussion in educational research. Test anxiety, known for disrupting cognitive processes and negatively impacting academic performance, has been studied extensively [7,8,29]. On the other hand, MSR is often associated with better academic outcomes, though its connection to test anxiety is less clear-cut [36,27]. Some research suggests that students who actively engage in MSR tend to experience lower levels of test anxiety, likely due to feeling more prepared and in control of their learning [40]. However, not all findings align with this view. In some cases, frequent use of MSR strategies has been linked to heightened stress, especially among students who are acutely aware of their academic challenges or feel external pressure to succeed [22]. Demographics further complicate this relationship [39,8]. Female students, for example, often report higher levels of test anxiety than males, even though they tend to use MSR strategies more effectively [23]. Despite these insights, the literature remains inconclusive on how generalizable these patterns are across different populations. This study aims to contribute to this ongoing debate.

2.2 Predictive Modeling of Anxiety Disorders

Machine learning has become a valuable tool for identifying individuals at risk of various anxiety disorders, enabling early and targeted interventions [12,42,1]. For example, Almadhor et al. [1] trained several models to predict anxiety levels, finding that Random Forest achieved the highest predictive accuracy. Similarly, Priya et al. [42] also applied machine learning to anxiety prediction, demonstrating strong performance in identifying negative cases. While these studies showcase the potential of predictive modeling, they tend to prioritize accuracy over fairness. Little attention is paid to whether predictions work consistently across diverse demographic groups, highlighting the need for research that considers both predictive accuracy and fairness to ensure that these tools serve all students effectively.

2.3 Fairness of Predictive Models in Education

Fairness in predictive modeling is an important issue in education [19,43,32,14]. Research shows that models optimized for accuracy may often perform worse for underrepresented groups [32,2,19]. Nonetheless there is a lack of consensus on the trade-off between fairness and predictive accuracy predictive even in the

general fair machine learning community [37]. Certain studies show that fairness and predictive accuracy can co-exist without a strict trade-off [19,20]. However, there are other studies that show that optimizing for fairness comes at cost to predictive accuracy [16,37,48]. These conflicting findings emphasize the need for further exploration, particularly in educational contexts where fairness is as important as predictive accuracy.

Furthermore, fairness evaluations in educational predictive models often neglect intersectionality [44,20,32]. To the best of our knowledge, only Gardner et al. [20] and Zambrano et al. [47] evaluate the fairness of predictive models in education along the intersection of multiple demographic attributes using the metric called AUC Gap. The AUC Gap performs well in highlighting intersectional subgroup disparity but it does not explicitly show whether there is a compounded (dis)advantage for a particular subgroup or not. The Residual Fairness Gap which we propose in Section 3.3, pinpoints intersectional subgroups with compounded (dis)advantages.

3 Methods

3.1 Data

We used survey data from the Motivated Strategies for Learning Questionnaire (MSLQ) [41], consisting of 81 items (i.e., survey questions) grouped into 15 subscales, collected over an 8-week period (April–June 2024) via Prolific with 672 consenting participants. Responses were rated on a 7-point Likert scale. Demographic data included race (54% White, 46% Non-White), migrant status (16% Migrants, 84% Non-Migrants), and sex (50.35% Male, 49.08% Female). In this study, we define migrants as individuals living in a country other than their birth country. For statistical power, race was categorized as White vs. Non-White, and sex analysis excluded the <1% who selected “Prefer not to say”. This study focuses on the *self-reported* Metacognitive Self-Regulation (MSR) and Test Anxiety (TA) subscales from the MSLQ, comprising twelve and five items respectively. The survey items have been rephrased for brevity, and we will use their aliases throughout this study (e.g., “*During class time I often miss important points because I’m thinking of other things*” is rephrased as `distracted_during_class`). See the supplementary sheet for the complete list [here](#).

Latent Test Anxiety Score Derivation: To represent test anxiety, we created a composite score from the five TA items. Internal consistency was verified with Cronbach’s Alpha ($\alpha = 0.80$) and McDonald’s Omega ($\omega = 0.81$), indicating strong reliability. Using exploratory factor analysis (EFA), we generated the composite score, confirming suitability with a Kaiser-Meyer-Olkin (KMO) test score of 0.80. The factor loadings (FLs) showed varying correlations of the five TA items with test anxiety, with `fear_of_failure` (FL= 0.81) being the most correlated. The factor scores were normalized between 0 and 1 to reflect increasing test anxiety.

3.2 Analyzing MSR, Test Anxiety, and Demographics

To investigate the relationship between test anxiety, metacognitive self-regulation, and student demographics, we conducted two statistical analyses. Firstly, we used Spearman’s rank correlation to analyze the relationship between twelve MSR behaviors and test anxiety. Despite some MSR features showing zero or near-zero correlation with test anxiety (e.g., `assess_concept_mastery`, $\rho = 0.00$), permutation importance (PMI) revealed that even these features could have predictive utility as shown in Table 2. Secondly, we used the Mann-Whitney U test to examine how race, sex, and migrant status affect test anxiety and MSR behaviors, calculating Cliff’s Delta (δ) to measure effect size and direction.

3.3 Predicting Test Anxiety

Predictive Models: To predict test anxiety, we selected five machine learning (ML) models commonly used in classification tasks in education: Support Vector Machine (SVM), Random Forest (RF), Logistic Regression (LR), XGBoost (XGB), and Multilayer Perceptron (MLP) [25,28,12,?]. These models have been effectively applied in similar contexts, such as predicting anxiety disorders and related mental health conditions [42,1].

Predictive Accuracy and Fairness Metrics: To evaluate the predictive accuracy of our models for identifying test anxiety, we used commonly applied metrics in machine learning—accuracy, F1 score, area under the receiver operating characteristic curve (AUC-ROC), and area under the precision-recall curve (AUC-PR) [19,43,32]. These metrics provide a comprehensive view, addressing the nuances of predicting test anxiety, where missing true cases or over-predicting can have serious consequences [5]. To assess fairness, we compared these metrics across demographic groups following established conventions [32,43,2], examining consistency and potential biases. Disparities in metrics between groups would indicate unfair outcomes, as the model may work better for some groups than others. Additionally, we investigate fairness along intersectional subgroups using our proposed metric

Proposed Metric-Residual Fairness Gap (RFG): RFG compares the *actual* predictive accuracy of an intersectional subgroup to the “*expected*” predictive accuracy of that subgroup, where the “expected” predictive accuracy is the overall average of the marginal predictive accuracies of its constituent groups. This approach reveals whether the intersection of demographic attributes introduces compounded effects—either positive or negative—on subgroup predictive accuracy. Given multiple sensitive attributes A, B, \dots, K , where a, b , etc., represent specific groups (e.g., a could represent “female” in sex and b could represent “Black” in race), RFG is defined as follows:

$$RFG_{a,b,\dots,k} = \text{Metric}_{a,b,\dots,k} - \frac{\text{Metric}_a + \text{Metric}_b + \dots + \text{Metric}_k}{n}$$
Where: **Metric** represents any predictive accuracy metric, such as precision or F1-score. **Metric_{a,b,...,k}** is the *actual* predictive accuracy of the intersectional subgroup (e.g., Black females) w.r.t the chosen metric, while **Metric_a, Metric_b, ..., Metric_k** are the marginal predictive accuracy scores of constituent groups w.r.t the chosen metric, and n is the number of sensitive attributes. The RFG score evaluates whether the intersectional subgroup performs better ($RFG > 0$, compounded advantage),

Table 1: *Ground truth* distribution showing the prevalence (or base rates) of test anxiety at different thresholds ($\tau = 0.4$, $\tau = 0.5$, and $\tau = 0.6$) across the overall population and different demographic groups in the dataset.

		Migrant Status			Race		Sex	
		Overall	Migrant	Non-Migrant	White	Non-White	Male	Female
Total Sample		672	107	565	363	309	338	334
Thresholds	$\tau = 0.4$	75.3%	80.4%	74.3%	75.8%	74.8%	74.6%	76.0%
	$\tau = 0.5$	62.9%	69.2%	61.8%	65.0%	60.5%	60.9%	65.0%
	$\tau = 0.6$	47.0%	50.5%	46.4%	49.0%	44.7%	45.3%	48.8%

worse ($RFG < 0$, compounded disadvantage), or as “expected” ($RFG = 0$, no compounded effects) relative to the overall average of its marginal groups.

Model Training and Evaluation: We built the test anxiety prediction models using the demographic variables and the metacognitive self-regulation items as input features, with binarized test anxiety scores as the target. Three binarization thresholds—0.4, 0.5, and 0.6—were used: 0.4 assumed false negatives to be costlier, 0.5 followed standard practice in classification, and 0.6 was derived from an ad-hoc ROC analysis for optimal sensitivity-specificity balance. We found that 0.4 overestimated test anxiety prevalence, while 0.6 underestimated it. This led us to choose 0.5 as the most balanced option as per Table 1. Nonetheless, we trained and tested all models using each threshold. In this paper, we will focus on 0.5 threshold, however, results for 0.4 and 0.6 thresholds are included in the supplementary materials for reference [here](#).

The five models were trained as follows. We split the dataset into an 80% training set and a 20% test set, using stratification to maintain the distribution of test anxiety. We determined optimal hyper-parameters through grid search with 5-fold cross-validation. After training, we evaluated the models on the test set by performing bootstrap sampling for 100 times to ensure robustness [18], ensuring that each bootstrap iteration contained the same *number* of observations as the original test set. For each iteration, we calculated the predictive accuracy and fairness metrics, then we averaged the results and calculated their standard deviations.

We assessed predictive accuracy using accuracy, F1 score, AUC-ROC, and AUC-PR. We evaluated fairness by computing differences in predictive accuracy across demographic groups and tested their significance using independent samples t-tests. Finally, to capture the compounded effect of bias across intersectional subgroups, we applied the RFG metric to bootstrap-averaged results for pairwise intersectional subgroups.

4 Results Discussion

4.1 RQ1: Dynamics of MSR, Test Anxiety, and Demographics

Recall that RQ1 explores the link between MSR, test anxiety, and their variation by demographics. As shown in Table 2, interestingly, we observed that out of

Table 2: Spearman’s rank correlation (ρ) between MSR features and test anxiety (TA), with PMI (permutation importance). Asterisks denote significance: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$ and is consistent throughout this study.

MSR	ρ	PMI
distracted_during_class	0.45***	0.78
formulate_guiding_questions	0.03	0.16
clarify_confusing_content	0.00	0.08
adjust_reading_strategy	-0.03	0.09
preview_course_material	-0.00	0.11
self_check_understanding	-0.01	0.12
adapt_study_methods	0.34***	0.35
mindless_class_reading	0.04	0.10
identify_learning_objectives	0.07	0.11
assess_concept_mastery	-0.00	0.11
set_study_goals	-0.05	0.11
review_unclear_notes	-0.01	0.07

the twelve MSR behaviors, only two—`distracted_during_class` ($\rho = 0.45$) and `adapt_study_methods` ($\rho = 0.34$)—are significantly ($p < 0.001$) correlated with test anxiety. However, it is not surprising to find that distraction during class was the most correlated MSR behaviour with test anxiety. We speculate that this could be due to the fact that students who are distracted during class may have missed important points that are crucial to understanding the study material, for example, due to mind wandering [17], thus feeling unprepared and anxious before tests [4]. Another interesting finding was that adapting study methods to fit course requirements or an instructor’s teaching style may not always reduce test anxiety, as previously found in studies such as [45]. Rather, we found that students who adapt their study methods may sometimes have high test anxiety. A probable reason for this could be that the increased pressure to adapt study methods or the frequency at which students keep changing study methods could destabilize their study routines [38]. This aligns with other research suggesting that frequent cognitive adjustments, especially when tied to metacognitive strategies, can increase test anxiety by adding to the mental load and making it harder for students to regulate their emotions effectively [21,22].

Across all demographic groups that we considered, i.e., race, sex, and migrant status, it was only in terms of sex that we observed significant difference in test anxiety. For instance, as shown in Table 3, the probability of a female reporting higher test anxiety was 10 percentage points greater than that of a male (approximately 55% vs. 45%). There are several existing studies that corroborates to this finding [39]. Focusing on the MSR behaviours that we found to be significantly correlated with test anxiety, i.e., `distracted_during_class` and `adapt_study_methods`, we did not find any significant difference across the various demographic groups except *race*. Specifically, we found that the White students are 15% more likely to be distracted during class compared to their Non-White counterparts. We conjecture that cultural differences can influence

Table 3: Mann-Whitney U test results for TA and MSR behaviours by Race (W: White, NW: Non-White), Sex (M: Male, F: Female), and Migrant Status (M: Migrant, NM: Non-Migrant). δ represents Cliff’s delta.

Variable	Race		Sex		Migrant Status	
	δ	Higher	δ	Higher	δ	Higher
Test Anxiety (TA)	0.03	W	-0.1*	F	-0.09	M
adapt_study_methods	0.04	W	-0.07	F	0.01	NM
adjust_reading_strategy	-0.21***	NW	-0.05	F	-0.02	M
assess_concept_mastery	-0.17***	NW	-0.11*	F	-0.07	M
clarify_confusing_content	-0.21***	NW	-0.13**	F	-0.03	M
distracted_during_class	0.15***	W	0.04	M	-0.07	M
formulate_guiding_questions	-0.29***	NW	-0.06	F	0.02	NM
identify_learning_objectives	-0.21***	NW	-0.11**	F	0.01	NM
mindless_class_reading	-0.19***	NW	0.03	M	-0.04	M
preview_course_material	-0.2***	NW	-0.1*	F	0.01	NM
review_unclear_notes	-0.18***	NW	-0.04	F	0.11	NM
self_check_understanding	-0.19***	NW	-0.08	F	0.05	NM
set_study_goals	-0.14**	NW	-0.16***	F	0.01	NM

Table 4: Average Predictive Accuracy of all models. The values are the the average \pm standard deviation. Boldened and red scores are the highest and least overall averages respectively

Model	AUC-ROC	AUC-PR	Accuracy	F1 Score	Overall Average
LR	0.66 \pm 0.05	0.76 \pm 0.05	0.67 \pm 0.04	0.75 \pm 0.03	0.71 \pm 0.04
MLP	0.54 \pm 0.05	0.67 \pm 0.05	0.61 \pm 0.04	0.71 \pm 0.03	0.63 \pm 0.04
RF	0.66 \pm 0.05	0.75 \pm 0.05	0.69 \pm 0.04	0.79 \pm 0.03	0.72 \pm 0.04
SVM	0.66 \pm 0.05	0.75 \pm 0.05	0.66 \pm 0.04	0.75 \pm 0.03	0.7 \pm 0.05
XGB	0.60 \pm 0.05	0.72 \pm 0.05	0.64 \pm 0.04	0.76 \pm 0.03	0.68 \pm 0.04

how students report being distracted. For example, White students might be more open about mentioning (even) minor distractions, while Non-White students, who are aware of biases and stereotypes [13,34], might downplay their distractions [46,26]. This is in line with studies showing that sociocultural norms affect emotional self-awareness and self-assessment [26]. Further future studies are needed to investigate this finding in detail.

4.2 RQ2: Trade-off Between Predictive Accuracy and Fairness

RQ2 examines whether the most accurate models for predicting test anxiety are also the fairest or involve trade-off. Firstly, in terms of predictive accuracy, no model consistently outperformed others across all metrics as shown in Table 4. However, averaging across metrics, the RF model performed best, while the deep

learning model (i.e., MLP) performed worst—contrary to prior studies where deep learning models excel [43]. This may be due to our dataset size, as neural networks often underperform on smaller datasets [6]. Nonetheless, similar to our results, several related studies have often found RF to outperform other models whenever such comparative analysis are done [28,35,9].

Table 5: Predictive accuracy disparities favor *historically* advantaged groups (Whites, Males, Non-Migrants) with negative values, and disadvantaged groups (Non-Whites, Females, Migrants) with positive values [2,32].

	Model	AUC-PR	AUC-ROC	Accuracy	F1 Score
Race	LR	0.02 ± 0.07 *	0.05 ± 0.08 ***	-0.07 ± 0.06 ***	-0.08 ± 0.06 ***
	MLP	0.06 ± 0.08 ***	0.13 ± 0.08 ***	0.02 ± 0.06 **	-0.01 ± 0.06
	RF	0.12 ± 0.06 ***	0.17 ± 0.07 ***	-0.02 ± 0.06	-0.02 ± 0.04 **
	SVM	0.09 ± 0.07 ***	0.11 ± 0.08 ***	-0.06 ± 0.06 ***	-0.07 ± 0.06 ***
	XGB	0.02 ± 0.08	0.08 ± 0.08 ***	0.03 ± 0.06 **	0.0 ± 0.06
Sex	LR	-0.02 ± 0.07	-0.12 ± 0.07 ***	-0.17 ± 0.06 ***	-0.13 ± 0.05 ***
	MLP	-0.09 ± 0.07 ***	-0.16 ± 0.08 ***	-0.11 ± 0.06 ***	-0.08 ± 0.06 ***
	RF	-0.03 ± 0.07 **	-0.1 ± 0.08 ***	0.0 ± 0.06	0.01 ± 0.04 *
	SVM	-0.06 ± 0.08 ***	-0.17 ± 0.08 ***	-0.14 ± 0.06 ***	-0.1 ± 0.06 ***
	XGB	-0.05 ± 0.06 ***	-0.15 ± 0.07 ***	0.03 ± 0.06 ***	0.04 ± 0.05 ***
MS	LR	0.05 ± 0.08 ***	-0.05 ± 0.1 **	0.02 ± 0.08	0.03 ± 0.07 ***
	MLP	0.04 ± 0.1 **	-0.11 ± 0.11 ***	-0.03 ± 0.08 *	0.0 ± 0.08
	RF	0.11 ± 0.08 ***	0.01 ± 0.1	0.0 ± 0.08	0.01 ± 0.06
	SVM	0.1 ± 0.08 ***	0.02 ± 0.1	0.04 ± 0.08 **	0.04 ± 0.07 ***
	XGB	0.06 ± 0.09 ***	-0.07 ± 0.1 ***	0.05 ± 0.08 ***	0.04 ± 0.06 ***

In terms of fairness, no model was consistently fair across all metrics and demographics as shown in Table 5. For instance, in terms of race, the LR model was *less* favorable to Whites by 2% (AUC-PR) and 5% (AUC-ROC) but *more* favorable to the same Whites by 7% (accuracy) and 8% (F1-score), supporting the fairness “*impossibility theorem*” which posits the mutual exclusivity of certain metrics [3,10,33].

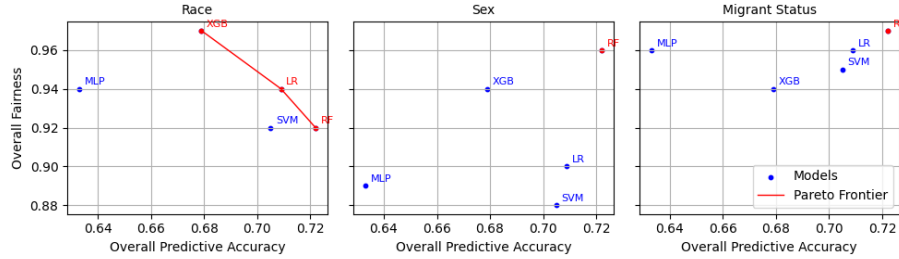


Fig. 1: Pareto frontier illustrating the trade-off between fairness and predictive accuracy. Red and blue points mark optimal and suboptimal models respectively.

We also observed a phenomenon that raises questions about how substantial some apparent cases of bias are. Specifically, we found that some disparities which we might be inclined to call bias are not statistically significant. For example, consider the fairness of RF in terms of race in Table 5. We observed that disparities in AUC-PR ($0.12 \pm 0.06^{***}$), AUC-ROC ($0.17 \pm 0.07^{***}$) and F1-score ($-0.02 \pm 0.04^{**}$) are significant. However, the disparity in accuracy (-0.02 ± 0.06) is not statistically significant. It is currently uncommon to use statistical significance testing to assess bias (but see [19]). Some differences may not be statistically significant, but that does not *necessarily* mean they are irrelevant. Better methods are needed to distinguish between disparities caused by random noise and those reflecting real bias. Without better methods, we risk missing real disparities or overreacting to random noise, which can lead to poor decisions about fairness.

Moving on to the crux of RQ2, we look at the trade-off between **overall fairness** and **overall predictive accuracy** across the three sensitive attributes: race, sex, and migrant status. Overall predictive accuracy of each model is operationalized as the mean predictive accuracy of that model across all evaluation metrics. Overall fairness of each model is operationalized using the absolute mean disparity in predictive accuracy between groups, normalized as $1 - |\text{mean disparity}|$. Using absolute values ensured that positive and negative differences did not cancel each other out, allowing us to capture the extent of unfairness regardless of its direction.

The results, as shown in Figure 1, highlight a clear trade-off between fairness and predictive accuracy. RF consistently achieved the highest predictive accuracy across all three attributes but did not always perform best in terms of fairness. For example, in the race analysis, XGB achieved the highest fairness but at the cost of slightly lower predictive accuracy. For sex and migrant status, RF was the only model on the Pareto frontier, meaning it offered the best balance between fairness and predictive accuracy, while other models like MLP and SVM were suboptimal, underperforming in both desiderata. Models below the frontier are less effective, as better-performing alternatives exist. Overall, the findings indicate that the most accurate models are not necessarily the most fair. Hence, improving fairness may often come at the expense of predictive accuracy. This finding partially contradicts studies [19,43] that reported no strict trade-off between predictive accuracy and fairness. However, numerous other studies, including ours, demonstrate that such a trade-off does exist [31,16,37].

4.3 RQ3: Intersectional Bias in Test Anxiety Models

This RQ aims to explore the compounded effect of the intersection demographic attributes. From the results in Table 6, we found that the intersection of demographic attributes can result in (1) compounded disadvantage, (2) compounded advantage, or (3) indifference, albeit, mostly in unexpected ways. For instance, *let us focus on the RF model* which we found to be the pareto optimal as per RQ2. From Table 5, for each demographic attribute in isolation, we found that the fairness of the RF model in terms of AUC-PR is: $0.12 \pm 0.06^{***}$ for race (advantage Non-Whites), $-0.03 \pm 0.07^{**}$ for sex (advantage Males), and $0.11 \pm 0.08^{***}$

Table 6: **RFG** for metrics: PR (AUC-PR), ROC (AUC-ROC), Acc (Accuracy), and F1 (F1 Score). Column initials indicate demographic intersections: NWM (Non-White Migrant), WM (White Migrant), NWNM (Non-White Non-Migrant), WNM (White Non-Migrant), NWF (Non-White Female), NWM (Non-White Male), WF (White Female), WM (White Male), FM (Female Migrant), MM (Male Migrant), FNM (Female Non-Migrant), MNM (Male Non-Migrant).

		Race-Migrant Status				Race-Sex				Sex-Migrant Status			
		NWM	WM	NWNM	WNM	NWF	NWM	WF	WM	FM	MM	FNM	MNM
LR	PR	0.15	-0.05	-0.01	0.02	-0.01	0.04	-0.01	-0.00	-0.04	0.13	0.02	-0.00
	ROC	0.18	-0.24	0.00	-0.01	-0.03	0.04	-0.12	0.00	-0.11	0.08	-0.01	0.02
	Acc	0.07	-0.09	-0.02	0.01	-0.11	0.07	0.04	-0.01	0.04	0.03	-0.03	0.04
	F1	0.06	-0.05	-0.03	0.02	-0.12	0.07	0.06	-0.02	0.02	0.04	-0.02	0.03
MLP	PR	0.13	-0.07	-0.00	0.03	-0.01	0.07	-0.02	-0.01	-0.01	0.13	0.01	0.02
	ROC	0.07	-0.27	0.02	-0.02	-0.01	0.07	-0.16	-0.01	-0.05	0.04	-0.04	0.04
	Acc	0.06	-0.17	-0.00	0.01	-0.06	0.07	-0.01	-0.02	-0.01	0.01	-0.03	0.03
	F1	0.06	-0.15	-0.02	0.02	-0.08	0.06	0.02	-0.03	-0.04	0.04	-0.02	0.01
RF	PR	0.11	-0.03	0.01	-0.00	0.00	0.07	-0.01	-0.04	0.02	0.07	0.00	0.00
	ROC	0.13	-0.23	0.02	-0.04	-0.04	0.09	-0.12	-0.06	-0.02	-0.03	-0.03	0.02
	Acc	0.07	-0.08	-0.01	0.01	-0.07	0.05	0.09	-0.07	-0.03	0.06	0.01	-0.01
	F1	0.05	-0.06	-0.02	0.01	-0.05	0.04	0.07	-0.05	-0.01	0.04	0.01	-0.01
SVM	PR	0.14	0.02	0.01	-0.02	-0.02	0.08	-0.04	-0.02	0.08	0.04	-0.02	0.02
	ROC	0.15	-0.11	0.01	-0.04	-0.05	0.09	-0.15	-0.03	0.06	-0.09	-0.05	0.05
	Acc	0.08	-0.08	-0.03	0.02	-0.07	0.04	0.02	0.01	0.04	0.05	-0.03	0.03
	F1	0.07	-0.04	-0.03	0.03	-0.09	0.04	0.04	-0.00	0.02	0.05	-0.02	0.02
XGB	PR	0.16	-0.13	-0.02	0.04	-0.02	0.06	0.03	-0.02	0.01	0.13	0.03	-0.00
	ROC	0.15	-0.36	-0.00	-0.01	-0.02	0.05	-0.10	0.01	-0.04	0.07	-0.02	0.03
	Acc	0.06	-0.05	-0.01	-0.00	-0.00	0.02	0.02	-0.03	0.10	-0.01	0.00	-0.00
	F1	0.05	-0.04	-0.02	0.01	-0.01	0.02	0.02	-0.03	0.06	0.01	0.01	-0.01

for migrant status (advantage migrants). With this in mind, one might expect the RF to have compounded advantage for Non-White Males, Non-White Migrants, and Male Migrants, for example. Similarly, one might expect compounded disadvantage for White Females, White Non-Migrants, and Female Non-Migrants. Our results in Table 6 sometimes agree with this hypothesis and at other times, disagree. For example, we observed that being Non-White Migrant resulted in a 11% improvement in predictive accuracy in terms of AUC-PR as compared to the average of the marginal AUC-PRs for Non-Whites and Migrants in isolation. In fact, across all models and all metrics, we observed a compounded advantage for Non-White Migrants. On the reverse, we expected compounded disadvantage for White Females. Truly, we observed a decline in predictive accuracy for White Females in terms of AUC-PR although it was a mealy 1% decline. Yet, in terms of some other metrics such AUC-ROC for the same RF model, we observed as high as a 12% decline in AUC-ROC for White Females relative the average of the marginal AUC-ROCs for Whites and Females in isolation. Nonetheless, we observed instances where there was actually an “unexpected” improvement in predictive accuracy for White Females for the same RF model in terms of

accuracy (9%) and F1-Score (7%). Overall, our results suggest that the combination of two supposed disadvantages as it were, may not necessarily result in compounded disadvantage as prior research [11] suggests and vice versa.

5 Concluding Discussions and Implications

Demographic Inequities in Test Anxiety and the Need for Targeted Interventions: Classroom distractions and inconsistent study methods have a significant negative impact on test anxiety. This is particularly concerning for White students, who are more prone to distractions, and for females, who experience higher anxiety levels compared to males. These patterns suggest potential inequities in how test anxiety manifests across demographics, emphasizing the need for targeted strategies to create focused learning environments and stabilize study routines. Moreover, the disproportionate effects of distractions and anxiety among specific groups imply that interventions should be tailored to address these demographic differences, such as designing inclusive classroom practices and offering gender-sensitive support programs.

Balancing Accuracy and Fairness in Predictive Models Through Thoughtful and Rigorous Evaluation: The trade-off between accuracy and fairness means that practitioners have to think carefully about what matters most for their specific goals. Is predictive accuracy the priority? Is fairness more important? Or is there a need to strike a balance between the two? For example, RF models were the most accurate overall, but they were not always the fairest—especially when looking at race, where XGB performed better in fairness but at a cost to predictive accuracy. Additionally, cases where disparities disappear after proper statistical testing highlight the importance of validating (un)fairness with care. Relying *solely* on raw differences can lead to unnecessary interventions that fail to address real issues. However, even with statistical testing, current methods are not fully reliable, and improving them is essential to ensure reliable and meaningful assessments of (un)fairness. Balancing accuracy and fairness requires thoughtful, evidence-based decisions, supported by transparent and robust evaluation methods.

Rethinking Intersectional Bias: Beyond Additive Assumptions in Predictive Fairness: Our findings indicate that the intersection of demographic attributes does not always lead to predictable compounded effects, challenging the assumption that combining disadvantages consistently exacerbates bias. This highlights the importance of moving beyond simple additive assumptions about fairness and adopting a more nuanced approach to understanding how intersectional attributes interact in predictive models. The variation in compounded effects further emphasizes the need for context-specific fairness evaluations, as the impact of intersectionality often depends on the metric or model, making tailored interventions essential to effectively addressing biases.

Limitation: We used binary demographic groupings due to the small dataset to ensure statistical power, but this limits the detection of nuanced group differences. Future work with a bigger dataset will enable finer categorization.

References

1. Almadhor, A., Abbas, S., Sampedro, G.A., Alsubai, S., Ojo, S., Al Hejaili, A., Strazovska, L.: Multi-class adaptive active learning for predicting student anxiety. *IEEE Access* (2024)
2. Baker, R.S., Hawn, A.: Algorithmic bias in education. *International Journal of Artificial Intelligence in Education* pp. 1–41 (2021)
3. Bell, A., Bynum, L., Drushchak, N., Zakharchenko, T., Rosenblatt, L., Stoyanovich, J.: The possibility of fairness: Revisiting the impossibility theorem in practice. In: *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. pp. 400–422 (2023)
4. Blankstein, K.R., Toner, B.B., Flett, G.L.: Test anxiety and the contents of consciousness: Thought-listing and endorsement measures. *Journal of Research in Personality* **23**(3), 269–286 (1989)
5. Bradford, A., Meyer, A.N., Khan, S., Giardina, T.D., Singh, H.: Diagnostic error in mental health: a review. *BMJ Quality & Safety* (2024)
6. Brigato, L., Iocchi, L.: A close look at deep learning with small data. In: *2020 25th international conference on pattern recognition (ICPR)*. pp. 2490–2497. *IEEE* (2021)
7. Cassady, J.C., Johnson, R.E.: Cognitive test anxiety and academic performance. *Contemporary educational psychology* **27**(2), 270–295 (2002)
8. Chapell, M.S., Blanding, Z.B., Silverstein, M.E., Takahashi, M., Newman, B., Gubi, A., McCann, N.: Test anxiety and academic performance in undergraduate and graduate students. *Journal of educational Psychology* **97**(2), 268 (2005)
9. Chen, J., Zhou, X., Yao, J., Tang, S.K.: Evaluation of student performance based on learning behavior with random forest model. In: *2024 13th International Conference on Educational and Information Technology (ICEIT)*. pp. 266–272. *IEEE* (2024)
10. Chouldechova, A.: Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big Data* **5**(2), 153–163 (2017)
11. Crenshaw, K.: Women of color at the center: Selections from the third national conference on women of color and the law: Mapping the margins: Intersectionality, identity politics, and violence against women of color. *Stanford Law Review* **43**(6), 1241–1279 (1991)
12. Daza, A., Saboya, N., Necochea-Chamorro, J.I., Ramos, K.Z., Valencia, Y.d.R.V.: Systematic review of machine learning techniques to predict anxiety and stress in college students. *Informatics in medicine unlocked* p. 101391 (2023)
13. Deckman, S.L.: Managing race and race-ing management: Teachers’ stories of race and classroom conflict. *Teachers College Record* **119**(11), 1–40 (2017)
14. Deho, O.B., Joksimovic, S., Li, J., Zhan, C., Liu, J., Liu, L.: Should learning analytics models include sensitive attributes? explaining the why. *IEEE Transactions on Learning Technologies* (2022)
15. Deho, O.B., Joksimovic, S., Liu, L., Li, J., Zhan, C., Liu, J.: Assessing the fairness of course success prediction models in the face of (un)equal demographic group distribution. In: *Proceedings of the Tenth ACM Conference on Learning @ Scale*. p. 48–58. *L@S ’23*, Association for Computing Machinery, New York, NY, USA (2023)
16. Deho, O.B., Zhan, C., Li, J., Liu, J., Liu, L., Duy Le, T.: How do the existing fairness metrics and unfairness mitigation algorithms contribute to ethical learning analytics? *British Journal of Educational Technology* (2022)
17. Desideri, L., Ottaviani, C., Cecchetto, C., Bonifacci, P.: Mind wandering, together with test anxiety and self-efficacy, predicts student’s academic self-concept but not

- reading comprehension skills. *British Journal of Educational Psychology* **89**(2), 307–323 (2019)
18. Efron, B., Rogosa, D., Tibshirani, R.: Resampling methods of estimation. In: Smelser, N.J., Baltes, P.B. (eds.) *International Encyclopedia of the Social & Behavioral Sciences*, pp. 13216–13220. Elsevier, New York, NY (2004)
 19. Gardner, J., Brooks, C., Baker, R.: Evaluating the fairness of predictive student models through slicing analysis. In: *Proceedings of the 9th International Conference on Learning Analytics & Knowledge*. p. 225–234. LAK19, Association for Computing Machinery, New York, NY, USA (2019)
 20. Gardner, J., Yu, R., Nguyen, Q., Brooks, C., Kizilcec, R.: Cross-institutional transfer learning for educational models: Implications for model performance, fairness, and equity. In: *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. pp. 1664–1684 (2023)
 21. Garnefski, N., Kraaij, V., Spinhoven, P.: Negative life events, cognitive emotion regulation and emotional problems. *Personality and Individual differences* **30**(8), 1311–1327 (2001)
 22. Ghribnavaz, S., Nouri, R., Moghadasin, M.: Relationship between metacognition believes and exam anxiety: Mediating role of cognitive emotion regulation. *Journal of Cognitive Psychology* **5**(4), 1–10 (2018)
 23. Gustems-Carnicer, J., Calderón, C., Calderón-Garrido, D.: Stress, coping strategies and academic achievement in teacher education students. *European Journal of Teacher Education* **42**(3), 375–390 (2019)
 24. Häuselmann, A., Custers, B.: Substantive fairness in the gdpr: Fairness elements for article 5.1 a gdpr. *Computer Law & Security Review* **52**, 105942 (2024)
 25. Hellas, A., Ihantola, P., Petersen, A., Ajanovski, V.V., Gutica, M., Hynninen, T., Knutas, A., Leinonen, J., Messom, C., Liao, S.N.: Predicting academic performance: a systematic literature review. In: *Proceedings companion of the 23rd annual ACM conference on innovation and technology in computer science education*. pp. 175–199 (2018)
 26. Huggins, C.F., Williams, J.H., Sato, W.: Cross-cultural differences in self-reported and behavioural emotional self-awareness between japan and the uk. *BMC Research Notes* **16**(1), 380 (2023)
 27. Huntley, C.D., Young, B., Tudur Smith, C., Fisher, P.L.: Metacognitive beliefs predict test anxiety and examination performance. In: *Frontiers in Education*. vol. 8, p. 1051304. Frontiers Media SA (2023)
 28. Hutt, S., Gardener, M., Kamantz, D., Duckworth, A.L., D’Mello, S.K.: Prospectively predicting 4-year college graduation from student applications. In: *Proceedings of the 8th International Conference on Learning Analytics and Knowledge*. pp. 280–289 (2018)
 29. Hutt, S., Ocumpaugh, J., Andres, J.M.A.L., Munshi, A., Bosch, N., Baker, R.S., Zhang, Y., Paquette, L., Slater, S., Biswas, G.: Who’s stopping you?—using micro-analysis to explore the impact of science anxiety on self-regulated learning operations. In: *Proceedings of the Annual Meeting of the Cognitive Science Society*. vol. 43 (2021)
 30. Hyseni Duraku, Z., Hoxha, L.: Self-esteem, study skills, self-concept, social support, psychological distress, and coping mechanism effects on test anxiety and academic performance. *Health psychology open* **5**(2), 2055102918799963 (2018)
 31. Jiang, W., Pardos, Z.A.: Towards equity and algorithmic fairness in student grade prediction. In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. pp. 608–617 (2021)

32. Kizilcec, R.F., Lee, H.: Algorithmic fairness in education. In: *The Ethics of Artificial Intelligence in Education*, pp. 174–202. Routledge (2022)
33. Kleinberg, J., Mullainathan, S., Raghavan, M.: Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807* (2016)
34. Kunesh, C.E., Noltemeyer, A.: Understanding disciplinary disproportionality: Stereotypes shape pre-service teachers’ beliefs about black boys’ behavior. *Urban Education* **54**(4), 471–498 (2019)
35. Lingjun, H., Levine, R.A., Fan, J., Beemer, J., Stronach, J.: Random forest as a predictive analytics alternative to regression in institutional research. *Practical Assessment, Research, and Evaluation* **23**(1), 1 (2019)
36. Mega, C., Ronconi, L., De Beni, R.: What makes a good student? how emotions, self-regulated learning, and motivation contribute to academic achievement. *Journal of educational psychology* **106**(1), 121 (2014)
37. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* **54**(6), 1–35 (2021)
38. Nicolas, T.M., Arambulo, R.: Test anxiety, readiness, and intervention strategies for enhancing board exam performance among psychology students. *The Quest: Journal of Multidisciplinary Research and Development* **2**(3) (2023)
39. Núñez-Peña, M.I., Suárez-Pellicioni, M., Bono, R.: Gender differences in test anxiety and their impact on higher education students’ academic achievement. *Procedia-Social and Behavioral Sciences* **228**, 154–160 (2016)
40. Onwunyili, F.C., Onwunyili, M.C.: Effect of self-regulated learning on test anxiety: Academic achievement and metacognition among secondary school students in anambra state. *South Eastern Journal of Research and Sustainable Development (SEJRSD)* **3**(2), 90–104 (2020)
41. Pintrich, P.: A manual for the use of the motivated strategies for learning questionnaire (mslq). National Center for Research to Improve Postsecondary Teaching and Learning (1991)
42. Priya, A., Garg, S., Tigga, N.P.: Predicting anxiety, depression and stress in modern life using machine learning algorithms. *Procedia Computer Science* **167**, 1258–1267 (2020)
43. Sha, L., Rakovic, M., Whitelock-Wainwright, A., Carroll, D., Yew, V.M., Gasevic, D., Chen, G.: Assessing algorithmic fairness in automatic classifiers of educational forum posts. In: *International Conference on Artificial Intelligence in Education*. pp. 381–394. Springer (2021)
44. Verger, M., Fan, C., Lallé, S., Bouchet, F., Luengo, V.: A comprehensive study on evaluating and mitigating algorithmic unfairness with the madd metric. *Journal of Educational Data Mining* **16**(1), 365–409 (2024)
45. Wittmaier, B.C.: Test anxiety and study habits. *The Journal of Educational Research* **65**(8), 352–354 (1972)
46. Xie, J.L., Roy, J.P., Chen, Z.: Cultural and individual differences in self-rating behavior: an extension and refinement of the cultural relativity hypothesis. *Journal of Organizational Behavior: The International Journal of Industrial, Occupational and Organizational Psychology and Behavior* **27**(3), 341–364 (2006)
47. Zambrano, A.F., Zhang, J., Baker, R.S.: Investigating algorithmic bias on bayesian knowledge tracing and carelessness detectors. In: *Proceedings of the 14th Learning Analytics and Knowledge Conference*. pp. 349–359 (2024)
48. Zhao, H., Gordon, G.J.: Inherent tradeoffs in learning fair representations. *Journal of Machine Learning Research* **23**(57), 1–26 (2022)