

# Auditing the Imputation Effect on Fairness of Predictive Analytics in Higher Education

Hadis Anahideh<sup>a</sup>, Nazanin Nezami<sup>a</sup>, Parian Haghghat<sup>a</sup>, Denisa Gàndara<sup>b</sup>

<sup>a</sup>*University of Illinois at Chicago, Department of Mechanical and Industrial Engineering, 842 W Taylor St, Chicago, 60607, IL, USA*

<sup>b</sup>*The University of Texas at Austin, Department of Educational Leadership and Policy, 1912 Speedway, Stop D5000, Austin, 78712, Texas, USA*

---

## Abstract

Colleges and universities use predictive analytics in a variety of ways to increase student success rates. Despite the potential for predictive analytics, two major barriers exist to their adoption in higher education: (a) the lack of democratization in deployment, and (b) the potential to exacerbate inequalities. Education researchers and policymakers encounter numerous challenges in deploying predictive modeling in practice. These challenges present in different steps of modeling including data preparation, model development, and evaluation. Nevertheless, each of these steps can introduce additional bias to the system if not appropriately performed. Most large-scale and nationally representative education data sets suffer from a significant number of incomplete responses from the research participants. Missing Values are the frequent latent causes behind many data analysis challenges. While many education-related studies addressed the challenges of missing data, little is known about the impact of handling missing values on the fairness of predictive outcomes in practice. In addition, assessing the fairness of predictive outcomes is an essential step in the development of equitable algorithms. In

particular, out-of-sample evaluation provides a less biased estimate of learning performance and fairness for real-time deployment.

In this paper, we set out to first assess the disparities in predictive modeling outcomes for college-student success, then investigate the impact of imputation techniques on the model performance and fairness using a commonly used set of metrics. We conduct a prospective evaluation to provide a less biased estimation of future performance and fairness than an evaluation of historical data. Our comprehensive analysis of a real large-scale education dataset reveals key insights on modeling disparities and how imputation techniques impact the fairness of the student-success predictive outcome under different testing scenarios. Our results indicate that imputation introduces bias if the testing set follows the historical distribution. However, if the injustice in society is addressed and consequently the upcoming batch of observations is equalized, the model would be less biased.

*Keywords:* machine learning, imputation, fairness evaluation, education

---

## 1. Introduction

Predictive analytics has become an increasingly hot topic in higher education. In particular, predictive analytics tools have been used to predict various measures of student success (e.g., course completion, retention, and degree attainment) by mapping the input set of attributes of individuals (e.g., the student's high school GPA and demographic features) with their outcomes (e.g., college credits accumulated) [1]. Campus officials have used these predictions to guide decisions surrounding college admissions and student-support interventions, such as providing more intensive advising to

certain students [1].

Despite the potential for predictive analytics, there is a critical disconnection between predictive analytics in higher education research and their accessibility of them in practice. Two major barriers to existing uses of predictive analytics in higher education that cause this disconnection are the lack of democratization in deployment and the potential to exacerbate inequalities.

First, education researchers and policymakers face many challenges in deploying predictive and statistical techniques in practice. These challenges present in different steps of modeling including data cleaning (e.g. imputation), identifying the most important attributes associated with success, selecting the correct predictive modeling technique, and calibrating the hyperparameters of the selected model. Nevertheless, each of these steps can introduce additional bias to the system if not appropriately performed [2]. Missing Values are the frequent latent causes behind many data analysis challenges. Most large-scale and nationally representative education data sets suffer from a significant number of incomplete responses from the research participants. While many education-related studies addressed the challenges of missing data, [3, 4, 5], little is known about the impact of handling missing values on the fairness of predictive outcomes in practice. To date, few works studied the impact of data preparation on the unfairness of the predictive outcome in a limited setting [6] or using merely a single notion of fairness metrics [7].

Second, predictive models rely on historical data and have the potential to exacerbate social inequalities [1, 8]. Over the last decade, researchers re-

alized that disregarding the consequences and especially the societal impact of algorithmic decision-making, might negatively impact individuals' lives. COMPAS, a criminal justice support tool, was found to be decidedly biased against Black people [9]. Colleges and universities have been using risk algorithms to evaluate their students. Recently, Markup investigated four major public universities and found that EAB's Navigate software is racially biased [10]. Ensuring fair and unbiased assessment, however, is complex and it requires education researchers and practitioners to undergo a comprehensive algorithm audit to ensure the technical correctness and social accountability of their algorithms.

It is imperative that predictive models are designed with careful attention to their potential social consequences. A wave of fair decision-making algorithms and more particularly fair machine learning models for prediction has been proposed in recent years years[11, 12]. Nevertheless, most of the proposed research either deals with inequality in the pre-processing or post-processing steps or considers a model-based in-processing approach. To take any of the aforementioned routes for bias mitigation, it is critical to audit the unfairness of the outcome of the predictive algorithms and identify the most severe unfairness issues to address.

Following these concerns, fairness audits of algorithmic decision systems have been pioneered in a variety of fields [13, 14]. The auditing process of unfairness detection of the model provides a comprehensive guideline for education researchers and officials to evaluate the inequalities of predictive modeling algorithms from different perspectives before deploying them in practice.

In this paper, we first study if predictive modeling techniques for student success show inequalities for or against a sample of marginalized communities. We use a real national-level education dataset to analyze the case of discrimination. We consider a wide range of Machine Learning models for student-success predictions. Then, we audit if prediction outcomes are discriminating against certain subgroups considering different notions of fairness to identify a potential bias in predictions. Furthermore, we investigate the impact of imputing the missing values using various techniques on model performance and fairness to key insights for educational practitioners for responsible ML pipelines. This study has the potential to significantly impact the practice of data-driven decision-making in higher education investigating the impact of a critical pre-processing step on predictive inequalities. In particular, how imputation impacts the performance and the fairness of a student-success prediction outcome.

In this paper, we present a prospective validation of predictive modeling on operational data as an important step in assessing the real-world performance of machine learning models. We conduct a prospective evaluation for Out-Of-Sample performance evaluation whereby a model learned from historical data is evaluated by observing its performance on new data. Prospective evaluation is likely to provide a less biased estimation of future performance than the evaluation of historical data. Ensuring that a model continues to perform well for integration into decision-making workflows and trustworthy operational impact, we simulate testing observations by adding subtle noise to a different group of predictors. To this end, we perturb the sensitive attribute *race*, non-sensitive attributes, and all predictors in three different

scenarios to demonstrate external generalization capacity to the community.

We predict the most common proxy attribute *graduation completion* concerning equal treatment to different demographic groups through different notions of fairness. The comprehensive study of the real large-scale dataset of ESL:2002 allows us to validate the performance of different ML techniques for predictive analytics in higher education in a real situation. To the best of our knowledge, none of the existing fair machine learning (ML) models have studied existing large-scale datasets for student success modeling. Most of the extant applications of fair ML demonstrate results using small datasets considering a limited number of attributes (e.g., [15, 16]) or in a specific context, such as law school admission [17, 18].

## 2. Bias in Education

*“Bias in, bias out”*. The first step toward auditing and addressing disparity in student success prediction is to understand and identify different sources of bias in the dataset. Most of the social data including education data are almost always biased since they inherently reflect historical biases and stereotypes [19]. Data collection and representation methods often introduce additional bias. Disregarding the societal impact of modeling the biased data, further exacerbates the discrimination in the predictive modeling outcome. The term bias refers to demographic disparities in the sampled data that compromise its representativeness [19, 20]. *Population bias* in the data prevents a model to be accurate for minorities [21]. Table 1 presents the

racial population bias in the ELS <sup>1</sup> dataset for students who attended four-year institutions. Therefore, students who identify as “White” are considered as the majority accounting for 66% of the observations, while “Multiracial”<sup>2</sup>, “Hispanic” and “Black” groups are underrepresented, making them minorities.

<b>Race</b>	<b>Percent of Population</b>
Asian	11.13%
Black	10.21%
Hispanic	8.13%
Multiracial	4.27%
White	65.83%

Table 1: ELS population bias

On the other hand, bias exists in the distribution of attributes’ values across different demographic groups, which is referred to as *behavioral bias*. Bias in the data has a direct impact on the algorithmic outcomes. One explanation is the strong correlation between sensitive attributes with other attributes and the response. Figure 1(a) shows the behavioral bias of the highest degree earned by different racial groups. The histogram indicates that among Black and Hispanic communities, degree attainment below the

---

<sup>1</sup>Education Longitudinal Study (ELS:2002) is a nationally representative, longitudinal study of 10th graders in 2002 and 12th graders in 2004. <https://nces.ed.gov/surveys/els2002/>

<sup>2</sup>We refer to students with two or more races as multiracial and denote it as MR

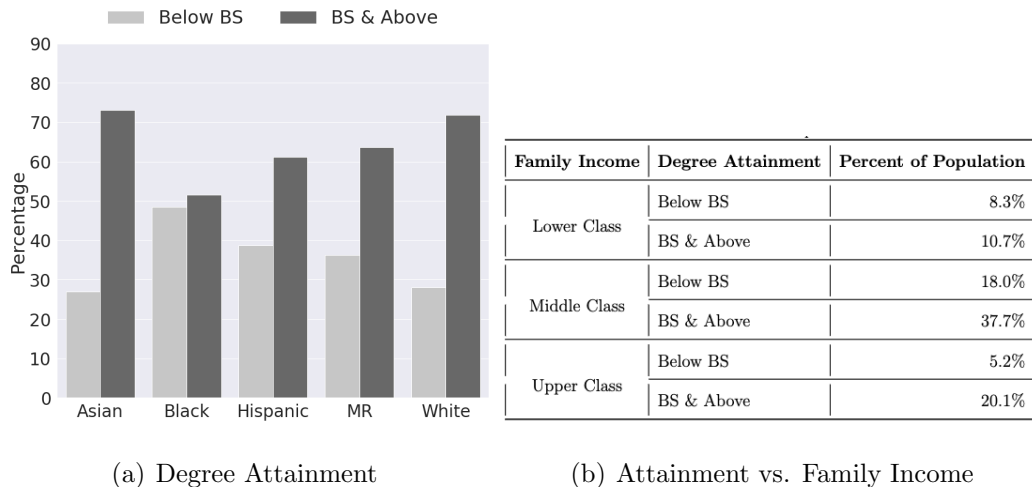


Figure 1: ELS behavioral bias

bachelor’s level is very frequent. We can observe another example of behavioral bias in Figure 1(b), where students from middle-class and low-income families have degree attainment lower than bachelor’s level. Therefore, using degree attainment as a student success indicator calls for careful action and thought to reduce the effects of both population and behavioral bias.

Transitioning toward a more comprehensive and specified model of disparities among population groups, yields a greater understanding of the key drivers of disparities among population groups. The key causes behind disparities among population groups that have been identified in the previous education research include, but are not limited to, social class[22, 23], race and ethnicity [22], gender [24, 25, 26], household characteristics(e.g. education of adults) [24], community characteristics (e.g. presence of schools) [24], and socioeconomic status[24, 27].

In this paper, we aim to investigate the impact of the data pre-processing step on disparities in the prediction outcome. Bias detection sheds light on



identifying the correct data pre-processing and imputation technique to alleviate the adverse impact of imputation on the performance and fairness of the model. To accomplish this end, we audit the unfairness of the predictive outcome before and after imputation using different scenarios that can occur in practice. We demonstrate how the obtained unfairness results differ and discuss the logic and implications behind each strategy. By examining bias in existing data and identifying the key characteristics of vulnerable populations, this paper illuminates how predictive models can produce discriminatory results if the bias is not addressed, and how we need to control the predictive outcome disparity.

To identify potential sources of behavioral bias, Figure 2(a), 2(b) and 2(c) illustrate racial disparities with respect to total credits, math/reading test scores, and GPA respectively. More specifically, Figure 2(a) shows that Black and Hispanic groups have lower median earned credits with their first and second quartile (50% of observations) plotted with lower values compared to others. Similarly, Figure 2(b) indicates that the student standardized combined math/reading test score has a lower median for Black and Hispanic groups. Moreover, the GPA score of vulnerable students (Black and Hispanic) have a lower median as shown in Figure 2(c). In addition, the size of each boxplot (from lower quartile to upper quartile) provides insight into the distribution of each group. For example, in Figure 2(a), the Black subgroup has a large box plot meaning that these students have very different outcomes in terms of total earned credits from very low to high values. However, the box plot for the White group of students indicates more similar credit outcomes, mainly distributed around the median value.

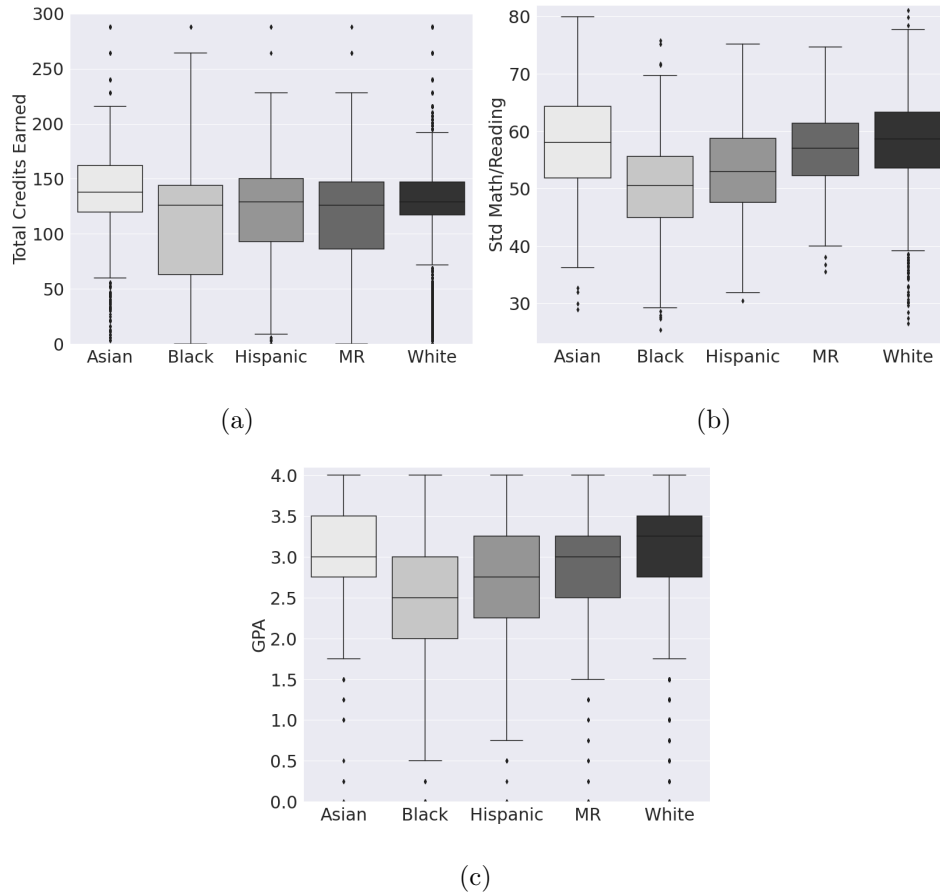


Figure 2: Boxplots of unprotected attributes  $\not\perp$  race

### 3. Fairness In Predictive Modeling

*Fairness-aware learning* has received considerable attention in the machine learning literature (fairness in ML) [28, 29]. More specifically, fairness in ML seeks to develop methodologies such that the predicted outcome becomes fair or non-discriminatory for individuals based on their protected attributes such as race and sex. The goal of improving fairness in learning problems can be achieved by intervention at pre-processing, in-processing

(algorithms), or post-processing strategies. Pre-processing strategies involve the fairness measure in the data preparation step to mitigate the potential bias in the input data and produce fair outcomes [30, 31, 32]. In-process approaches [33, 34, 35] incorporate fairness in the design of the algorithm to generate a fair outcome. Post-process methods [36, 30, 37], manipulate the outcome of the algorithm to mitigate the unfairness of the outcome for the decision-making process.

There are various definitions for fairness in the literature [38, 20, 39, 2, 40]. The fairness definitions fall into different categories including Statistical Parity [41], Equalized Odds [41], Predictive Equality [42] and Equal Opportunity [43]. Table 2 demonstrates the mathematical definitions of each of these common metrics [44]. Let  $S = \{0, 1\}$  be a binary sensitive attribute, in a binary classification setting let  $Y = \{0, 1\}$  be the true label, and  $\hat{Y} = \{0, 1\}$  be the predicted class label. Most of the fairness notions are derived based on conditional probabilities using these variables to reveal the inequalities of the predictive model. Evaluating the fairness of algorithmic predictions requires a notion of fairness, which can be difficult to choose in practice. Different metrics have been leveraged regarding different contexts, business necessities, and regulations. A predictive modeling outcome might have inequalities under one notion of fairness and might not have any under others.

In the education context, to give some examples, a) Demographic (statistical) parity is referred to as the *discrepancy of the predicted highest level of degree (success) across different demographic groups of students*, and b) Equal opportunity indicates the *discrepancy of the predicted highest level of a degree across different demographic groups of students, given their success*,

Fairness Notion	Formulation
Statistical Parity ( <b>SP</b> )	$P(\hat{Y} = 1 S = 1) - P(\hat{Y} = 1 S = 0)$
Equalized Odds ( <b>EO</b> )	$P(\hat{Y} = 1 Y = y, S = 1) - P(\hat{Y} = 1 Y = y, S = 0), \forall y \in \{0, 1\}$
Equal Opportunity ( <b>EO</b> P)	$P(\hat{Y} = 1 Y = 1, S = 1) - P(\hat{Y} = 1 Y = 1, S = 0)$
Predictive Equality ( <b>PE</b> )	$P(\hat{Y} = 1 Y = 0, S = 1) - P(\hat{Y} = 1 Y = 0, S = 0)$

Table 2: Common Fairness Definitions

is 1. In this paper, we use a binary classification setting but across multilevel racial and gender population subgroups ( $S$  is not necessarily binary).

We extend the fairness metrics described in Table 2 for non-binary sensitive attributes, by considering *one-versus-rest* approach for unfairness calculation. More specifically, to calculate the unfairness gaps, we consider each subgroup as  $S = 1$  and compare it against the rest  $S = 0$  (i.e. all other subgroups), one at a time. In this paper, we mainly focus on racial disparities, however, our proposed approach for auditing fairness and investigating the imputation impact can be extended to use other sensitive attributes. For example, the decision maker can use "gender" as a sensitive attribute.

#### 4. Fairness Audits

Notwithstanding the awareness of biases and unfairness in machine learning, the actual challenges of ML practitioners have been discussed in a few previous research [45, 46] with the focus on specific contexts, such as predictive policing [9] and child mistreatment detection [47]. ML practitioners often struggle to apply existing auditing and de-biasing methods in their contexts [46]. The concept of auditing algorithms and ethics-based auditing in

various contexts has lately reached its pinnacle. [48, 49, 50, 13, 48]. The final goal of the fairness auditing process is to determine whether the ML model’s results are fair. As a result, the auditing process aids in determining the appropriate actions to take, the best bias mitigation method to employ, and the most suitable technique to use throughout the ML development pipeline [49].

A designated user of predictive modeling in higher education needs support to audit the ML model performance and inequalities before adopting and deploying it in practice. To address the education practitioners and policymakers on assessing the inequalities of the predictive outcome, in this paper, we audit the unfairness of ML models for student success prediction using major notions of fairness to identify bias in algorithms using different metrics, Table 2. We also audit unfairness to ensure an ethical pre-processing approach. We audit a wide range of fairness metrics and conduct a comprehensive analysis of the performance of different ML models and their inequalities across different racial and gender subgroups throughout the data preparation (imputation) and the model training steps using the ELS dataset.

## **5. Success Prediction**

Before moving to the ML pipeline, we first discuss the prediction problem of interest. In this paper, we specifically focus on predicting the academic success of students in higher education. Student-success prediction is critical for institution performance evaluation, college admission, intervention policy design, and many more use cases in higher education [16, 51].

Quantifying student success is a very complex topic since the true quality

of any candidate is hidden and there is limited information available. There are proxy attributes such as first-year GPA or graduation completion that are typically used as the measure of success. In this work, we are primarily interested in studying the prediction of the highest level of degree (classification problem) using the ELS:2002 dataset.

Numerous factors can affect student success [26]. Thus, identifying the most informative and significant subset of potential variables is a critical task in predictive modeling [52]. To select a proper subset of attributes, we conducted a thorough literature search and combine it with the domain expert knowledge. These factors include, but are not limited to, academic performance (SAT scores, GPA) [53], student demographic attributes (e.g. race, gender) [26, 54], socio-economic status [24, 55], environmental factors, and extra(out of school) activities.

Incorporating protected attributes in the modeling procedure has raised concerns in the fair-ML domain [2]. The predictive outcome depends on the information available to the model and the specific algorithm used. A model may leverage any feature associated with the outcome, and common measures of model performance and fairness will be essentially unaffected. In contrast, in some cases, the inclusion of unprotected attributes may adversely affect the performance and fairness of a model due to a latent correlation with other protected attributes. In this paper, we audit the unfairness of the model and the impact of imputation incorporating the sensitive attribute as the determinant. Decision Tree [56], Random Forest [57], K-Nearest Neighbor [58] [59], Logistic Regression [60], and SVM [61] are among the well-known ML models in higher education. Table 3, represents the list of variables in

this study, and their corresponding missing value percentages.

Variables	% Missing	Variables	% Missing
S-T relationship	33.32	Std Math/Reading	0.02
F3-loan-owed	25.33	F3_Employment	0
%white teacher	23.69	F3_Highest level of education	0
%Black teacher	19.85	High school attendance	0
%Hispanic teacher	17.72	Family Composition	0
TV/video(h/day)	14.87	Race_Hispanic, race specified	0
Work(h/week)	12.06	F3_Separated no partner	0
F2_College entrance	9.75	F3_Never Married w partner	0
Generation	7.06	F3_Never Married no partner	0
F3_GPA(first attended)	6.79	F3_Married	0
F3_GPA(first year)	6.78	F3_Divorced/Widowed w partner	0
F1_TV/video(h/day)	6.76	F3_Divorced/Widowed no partner	0
F1_units in math	6.13	Race_White	0
Athletic level	5.39	Race_More than one race	0
F1_frequency of computer use	4.27	Race_Hispanic, no race specified	0
Credits (total)	4.04	School Urbanicity	0
F1_Std Math	3.64	Race_Black or African Amer	0
Credits (first year)	3.48	Race_Asian, Hawaii/Pac. Isl	0
F3_GPA (all)	3.33	Race_Amer. Indian/Alaska	0
F3_Credits_Math	3.01	Gender_Male	0
F3_Credits_Science	2.90	Gender_Female	0
F1_Work(h/week)	2.77	Parents education	0
Homework(h/week)	1.85	Income	0
Number of school activities	0.84	F1_Drop out	0
English	0.02	F3_Separated w partner	0

Table 3: List of Variables



### 5.1. Missing Values and Imputation

*Missing Values* are the frequent latent causes behind many data analysis challenges, from modeling to prediction, and from accuracy to fairness for protected (sensitive) subgroups. Therefore, *handling Missing values* is a complicated problem that requires careful consideration in education research [3, 4]. In this regard, different imputation techniques have been proposed in the literature and the effectiveness of each methodology on various applications has been studied. *Mean Imputation*, *Multiple Imputation*, and other clustering-based imputation strategies such as *KNN-imputation*, are among the well-known techniques in handling missing values which we will briefly describe here.

Most large-scale and nationally representative education data sets, (e.g., ELS) suffer from a significant number of incomplete responses from the research participants. While features that contain more than 75% missing or unknown values are not usually informative, most features suffer from less than 25% missing values and are worth keeping. Removing all observations with missing values induces significant information loss in success prediction.

**Simple Imputation** is one the most basic imputation strategies. The process involves replacing the missing variable of observation with the mean (or median) of the observations with available values for the same variable. The mean imputation method is known to decrease the standard error of the mean. This fact exacerbates the risk of failing to capture the reality through statistical tests [62, 3]. **Multiple Imputation (MI)** [63] is a more advanced imputation strategy that aims to estimate the natural variation in the data by performing several missing data imputations. In fact, MI produces a set of

estimates through various imputed datasets and combines them into a single set of estimates by averaging across different values. The standard errors of parameter estimates produced with this the method has shown to be unbiased [63]. **KNN Imputation** is a non-parametric imputation strategy, which has been shown to be successful for different contexts [64]. KNN imputer replaces each sample’s missing values with the mean value from  $K$  nearest neighbors found in the dataset. In fact, two samples are considered close neighbors if the features that neither is missing are close. KNN imputation is able to capture structure in the dataset while the underlying data distribution is unknown [65]. To the best of our knowledge, KNN imputation has not been used in the education context.

Overall, ignoring missing data can not be an effective approach of handle missing values, and more importantly, can result in predictive disparity for minorities. While many education-related studies addressed the challenges of missing data, as discussed, little is known about the impact of applying different imputation techniques on the fairness outcome of the final model. In this project, we aim to address this gap by considering the three above-mentioned imputation strategies. To the best of our knowledge, none of the prior works in the education domain worked on ensuring fairness while imputing missing values in pre-processing steps.

## 6. Experiments

**Data Preparation.** As previously stated, we use the ELS dataset in this study to audit the fairness of ML models in the development pipeline for predicting student success. The ELS dataset includes many categorical variables. Therefore, we begin by creating appropriate labeling and converting

categorical attributes to numeric ones (dummy variable) following the NCES dataset documentation<sup>3</sup>). Next, we perform a transformation on the considered response variable *highest level of degree* to construct a binary classification problem. That is, we label students with a college degree (BS degree and higher) as the favorable outcome (label=1), and others as the Unfavorable outcome (label=0). Note that the dataset is filtered based on the institution type to only include students who attended four-year postsecondary institutions. We consider “race” as the sensitive attribute for fairness evaluation, which includes five groups of White, Black, Hispanic, Asian, and Multi-racial(MR) students. A Data cleaning is performed to identify and rename the missing values (based on the documentation) and remove the observations that have many missing attributes (> 75% of the attributes are missing).

**Imputation.** The final and significantly important task is then to handle the remaining missing values in the dataset. We consider different imputation techniques; Simple Imputation **SI**, Multiple Imputation **MI**, and KNN Imputation **KNN-I**. We consider a baseline where we remove observations with missing attributes and refer to it as **Remove-NA**.

**Model Training** procedure follows the data preparation step as we obtain clean-format datasets. We aim to analyze the performance of different Machine learning (ML) models under each imputation technique and data split scenario, to audit the inequalities in the prediction outcome. We consider Decision Tree **DT**, Random Forest **RF**, Support Vector Classifier **SVC**, and

---

<sup>3</sup>[https://nces.ed.gov/surveys/els2002/avail\\_data.asp](https://nces.ed.gov/surveys/els2002/avail_data.asp)

Logistic Regression **Log** ML models.

### 6.1. Prospective Evaluation

To investigate the impact of imputation and data representation on the accuracy and fairness of the outcome of the predictive modeling, we simulate various scenarios to provide a trustworthy evaluation of models and their generalization capacity for real-time deployment. Table 4 summarizes our considered scenarios to evaluate imputation impact on the model unfairness. **RNA.rnd** indicates the common practice data pre-processing step in the machine learning literature where all rows with missing values are removed before a train/test split. Similarly, in **RNA.str** we remove the missing values while performing stratification on the “race” and “response” (highest level of degree) variables to ensure that students from each racial subgroup with different success outcomes (response) are well-represented in both training and testing datasets.

Scenario	Missing Values	Train Test Data
RNA.rnd	Removed	80:20 split
RNA.str	Removed	80:20 split, stratified on race and target variable
Imp.rnd	Imputed	80:20 split
Imp.str	Imputed	80:20 split, stratified on race and target variable
Imp.prop	Imputed	80:20 split, dataset separated based on racial groups, where training/testing data is proportional to the fraction of observations within each group
Imp.prop.frac	Imputed	Similar to Imp.prop scenario but only 75% of the testing data is chosen (randomly) for testing
Imp.prop.frac.perturb	Imputed	Similar to Imp.prop.frac where the dataset is also simulated by performing three different perturbations on race, nonsensitive, and all attributes

Table 4: Evaluation Scenarios

In **Imp.rnd**, we split the entire dataset into train/test sets, perform an imputation technique on the training set, and transfer the trained imputer on the testing set to replace missing values. In **Imp.str** we perform imputation similar to **Imp.rnd**, however, we consider stratification on the “race” and “response” variables to generate representative train and test datasets. In another attempt to maintain the distribution of different population subgroups, the **Imp.prop** scenario by fixing the fraction of observation from each racial group in train/test splits to the fractions in the entire dataset. To simulate a real-time scenario, we randomly select a smaller subset of the remaining observations for testing in **Imp.prop.frac** since there is no guarantee that future observations to maintain the same distribution as training data. Lastly, in **Imp.prop.frac.perturb**, a modified testing dataset is generated through perturbation of different groups of attributes to simulate the representation of different racial groups and their associated attribute values for future observations. That is, we perturb “race”, “non-sensitive” attributes, and “all attributes”, respectively.

## 6.2. Results

In this section, we summarize our findings in three key discussions. First, we compare the performance and unfairness of the considered ML models using different imputation strategies. Next, we compare random and stratified imputed data against **Remove-NA** scenarios. We then simulate different perturbation scenarios to study the impact of imputation on the unfairness while considering unknown distribution for future observations (real-time scenario). Lastly, we compare different pre-processing approaches and real-time scenarios to offer general guidelines for education practitioners.

Figure 3 reveals the impact of utilizing different imputation techniques on the unfairness of different ML models based on the Statistical Parity (**SP**) notion. We can observe that the unfairness has increased for Black and Hispanic groups of students after imputation. It is worth mentioning that imputation significantly reduces the variance of unfairness across different racial groups as it involves more observations. The unfavorable unfairness hike is still true for all ML models, if we solely change the imputation strategy. As a result, even while imputation has a considerable impact on the level of unfairness (raising the unfairness mean and decreasing the variance), the imputation techniques evaluated in this study are barely different from one another. Additionally, the impact of various imputation methods on unfairness exhibits a consistent pattern based on other fairness notions presented in § 8.

Figure 4 presents a comparison between **RNA-str** and **RNA-rnd** with imputing the missing values with **Imp-str** and without **Imp-rnd** stratification. Although imputation improves the prediction accuracy, it hurts the unfairness gap for Black and Hispanic groups. For example, the average unfairness based on **SP** increases in magnitude for both Black and Hispanic students after imputation (**Imp-str** and **Imp-rnd**) compared to Remove-NA scenarios (**RNA-str** and **RNA-rnd**). One justification behind this observation can be related to the size of the testing set, which significantly increases after imputation, Table 6.2. The average unfairness gaps increase as we include more observations from the minority (underprivileged) group of students via imputation.

Figure 5 presents the results of considered unfairness metrics for differ-

Scenario	Race	Pop. $X_{train}$	Pop. $X_{test}$	$Y = 1$	$\hat{Y} = 1$
<b>RNA.rnd</b>	Asian	95	26	20	22
	Black	75	19	13	15
	Hispanic	75	17	13	14
	Multiracial	32	9	7	8
	White	712	175	128	145
<b>Imp.rnd</b>	Asian	494	123	90	108
	Black	452	114	62	77
	Hispanic	357	94	58	67
	Multiracial	190	47	30	35
	White	2923	727	526	598
<b>Imp.prop</b>	Asian	494	123	90	100
	Black	452	114	58	66
	Hispanic	361	90	55	65
	Multiracial	190	47	30	36
	White	2920	730	525	614

Table 5: Population size of demographic groups

ent evaluation scenarios. Recall that the perturbation scenarios simulate real-time testing data, which may or may not follow the original distribution in the original dataset. That is, we merely perturb the “race” in **Imp.prop.frac.perturb-race**, all non-sensitive attributes (all attributes except race and gender) in **Imp.prop.frac.perturb-nonsensitive**, and all attributes in **Imp.prop.frac.perturb-all**. First, we notice that perturbation based on race decreases unfairness, compared to **RNA.rnd** and **Imp.rnd**, while increasing the accuracy. This indicates that perturbing the sensitive

attribute in the testing data first equalizes the representation of all groups. As a result, it decorrelates the sensitive and non-sensitive attributes inducing unfairness reduction and accuracy improvement. In addition, perturbing the non-sensitive attributes decreases unfairness for Black and Hispanic racial groups while decreasing the model accuracy. Note that unfairness reduction is less significant than perturbing race only. It is possible that only the correlation between sensitive and non-sensitive attributes has been decreased, leaving the representation bias (population bias) unchanged when non-sensitive attributes are perturbed. A similar pattern can be seen in perturbing all attributes (the most realistic case), which will result in a less-accurate and more fair model compared with the **RNA-rnd** and **Imp-rnd** scenarios. Lastly, perturbing all of the attributes of the not-imputed dataset, **RNA.rnd.perturbed-all**, would reduce the unfairness for vulnerable groups, however, it adversely affects privileged groups. That is because the test set for **RNA.rnd.perturbed-all** scenario suffers from representation bias.

The perturbation results shown in Figure 5 reflect scenarios where social injustice has been mitigated through effective interventions such as increasing access to more educational resources (e.g. computers and scholarships) for vulnerable population groups to increase their representation in the applicant pool and boost their chances of success (attain bachelor degree). In such an ideal scenario, the model will will treat every group equally. Note that perturb-all and non-sensitive scenarios have low accuracy since distribution of the input attribute space differs from the distribution used to build the model.



In Figure 6, we compare the estimated unfairness gaps for perturbation, regular and modified imputation, and compared them against Remove-NA scenarios. First, we note that the unfairness gaps are larger for the imputation scenarios (**Imp-rnd** and **Imp-prop** and **Imp-prop-frac**) compared to the **RNA-rnd** scenarios. Furthermore, perturbation scenarios show that the high unfairness gaps in imputation scenarios might have been overestimated since the distribution of future observation may differ dramatically from the past if interventions are implemented in society to address the disparity. By perturbing the data distribution in the test set, we eliminated representation bias and the correlation between sensitive and unprotected attributes. First, the testing distribution is different than training, and, second, the test set observations exhibit identical characteristics, meaning there is no difference between population subgroups. Note that the model performs similarly poorly (testing accuracy) for all subgroups in this scenario as the testing distribution is different than training (i.e., “no free lunch theorem”).

As shown in Table 2, Statistical Parity (**SP**) is a fairness metric that compares positive (or success) prediction outcome ( $\hat{Y} = 1$ ) across different racial groups without considering their true  $Y$  label (real outcome). Figure 4 shows that the unfairness for Black and Hispanic students increases after imputation based on (**SP**), which suggests that fewer students are being predicted as positive (successful) in comparison to other racial groups. On the other hand, Figures 5 and 6 indicate that if we perturb all the attributes in the testing data (i.e., enforcing effective interventions), Black and Hispanic students tend to be predicted as positive (successful), thereby reducing the unfairness gap in prediction.

Predictive Equality (**PE**) focuses mainly on unsuccessful students who are incorrectly predicted as successful given their racial subgroup. This type of unfairness could lead to considerably unfavorable results in higher education, where policymakers fail to recognize those in need. Figure 4 reveals that unfairness for Black and Hispanic groups increases after imputation, meaning that more unsuccessful students are incorrectly predicted as positive (successful). Results in Figures 5 and 6 show that imputation without consideration of intervention (changing the distribution with perturbation) leads to an almost similar level of unfairness for Black and Hispanic students. However, the unfairness decreases with perturbation/simulated scenario leading to insignificantly small incorrect classification of unsuccessful students from these groups.

Equal Opportunity (**EoP**) also emphasizes the positive prediction outcome and examines how the model correctly classifies successful students given their racial subgroup. Figure 4 presents that unfairness for Black and Hispanic students does not change considerably with and without imputation since the model correctly classifies the positive class (successful students) as the model accuracy level remains roughly the same. Moreover, Figures 5 and 6 further demonstrate that the degree of unfairness based on **EoP** for the perturbation scenario may or may not decrease for Black. The justification lies behind the impact of perturbation on the accuracy level of the model. Changing the distribution reduced the probability of correct classification of successful students (True positives), meaning that they are less correctly classified as successful. Therefore, simulating the race reduces unfairness based on **EoP** through decorrelating sensitive and non-sensitive attributes

and improving the generalization power of the model (accuracy level). However, perturbing all the attributes does not necessarily decrease the unfairness gaps for the unprivileged group (Black and Hispanic students) since changing the distribution of all attributes results in less accurate predictions.

Equalized Odds (**EO**) measures the True Positive and False Positive rates across different racial groups. It merely emphasizes the positive prediction outcome  $\hat{y} = 1$ , which is *>BS degree attainment* conditioned on the true response value and race. Figure 4 shows that unfairness for Black and Hispanic students increases after imputation based on (**SP**) which indicates that more students are predicted as positive (successful) compared to other racial groups, which can lead to adverse outcomes if the prediction is a false positive one. On the other hand, Figures 5 and 6 indicate that if we perturb all of the attributes in the testing data (real-time case) the Black and Hispanic groups tend to be less predicted as positive (successful) reducing the unfairness and notifying the decision-makers of the required policy changes needed to be considered in future.

## 7. Conclusion

Education policymakers face numerous challenges, including data preparation when utilizing predictive modeling to aid in decision-making processes. In addition, deploying machine learning on historical data would introduce bias to future predictions. One primary source of bias originates from the representativeness of vulnerable population subgroups in the historical data. In this paper, we assessed the impact of handling missing values with imputation on the prediction unfairness for the large-scale national ELS dataset.

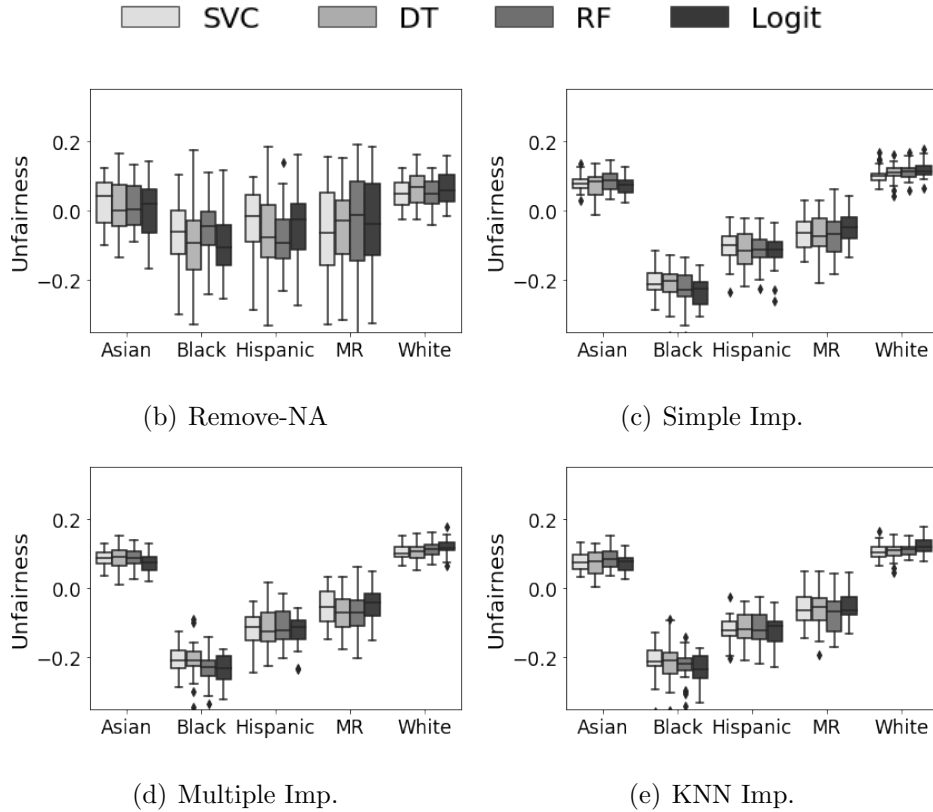


Figure 3: Impact of Imputation - Statistical Parity

First, we compared the standard procedure of removing all rows with missing values with the regular train-test split process for imputed missing values. In comparison to the Remove-NA scenario, imputation increases the average while reducing the variance of the unfairness gap for vulnerable groups of students. On the other hand, imputation increases the prediction accuracy of different ML models, indicating a trade-off with fairness. Moreover, we observed that the choice of imputation strategy has an insignificant impact on the unfairness gaps for different population subgroups.

To evaluate the unfairness of the ML model outcome, we designed var-

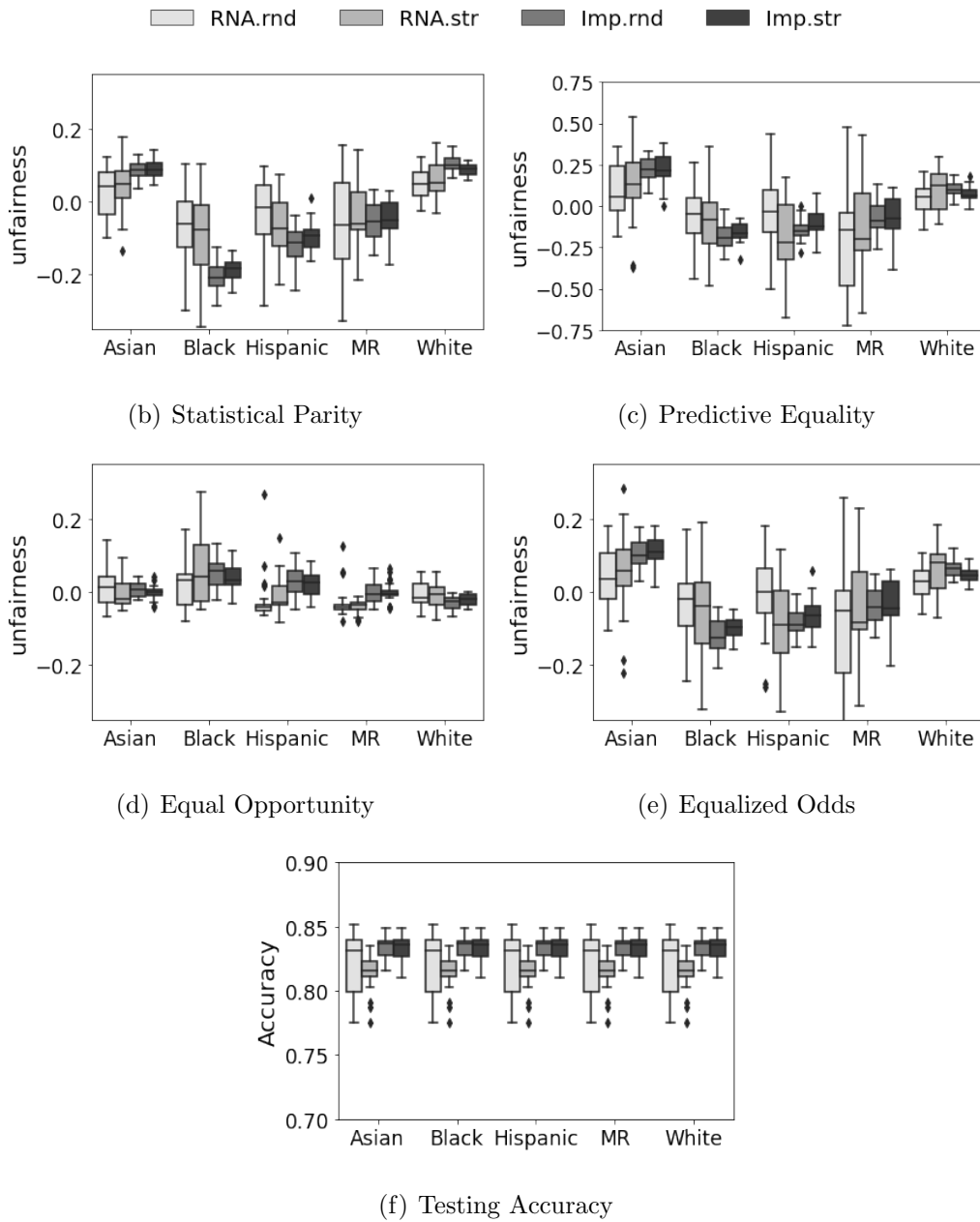


Figure 4: Remove-NA vs. Imputation (SVC)

ious train-test splitting scenarios and altered the distribution of future observations by simulating the testing data. Our analysis indicated that while

imputation tends to increase the average unfairness when using the standard training-testing split approach, it assists with attaining fair models when the testing distribution changes for future observations (as a result of effective disparity reduction interventions in society). In other words, if access to educational resources, parental education, community training, and general interventions to support vulnerable students are improved, the upcoming batch of students would have a different distribution. Thereby, the model trained on the imputed training data would be fair enough.

## References

- [1] M. Ekowo, I. Palmer, The promise and peril of predictive analytics in higher education: A landscape analysis., *New America* (2016).
- [2] S. Barocas, A. D. Selbst, Big data’s disparate impact, *Calif. L. Rev.* 104 (2016) 671.
- [3] J. R. Cheema, A review of missing data handling methods in education research, *Review of Educational Research* 84 (4) (2014) 487–508.
- [4] C. A. Manly, R. S. Wells, Reporting the use of multiple imputation for missing data in higher education research, *Research in Higher Education* 56 (4) (2015) 397–409.
- [5] S. K. Kwak, J. H. Kim, Statistical data preparation: management of missing values and outliers, *Korean journal of anesthesiology* 70 (4) (2017) 407.
- [6] I. Valentim, N. Lourenço, N. Antunes, The impact of data preparation on the fairness of software systems, in: *2019 IEEE 30th International Symposium on Software Reliability Engineering (ISSRE)*, IEEE, 2019, pp. 391–401.
- [7] M.-P. Fernando, F. Cèsar, N. David, H.-O. José, Missing the missing values: The ugly duckling of fairness in machine learning, *Journal of Intelligent Systems* (2021).
- [8] R. F. Kizilcec, H. Lee, Algorithmic fairness in education, *arXiv preprint arXiv:2007.05443* (2020).

- [9] J. Angwin, J. Larson, S. Mattu, L. Kirchner, Machine bias: Risk assessments in criminal sentencing, ProPublica (2016).
- [10] T. Feathers, Major universities are using race as a “high impact predictor” of student success, [themarkup.org/news/2021/03/02/major-universities-are-using-race-as-a-high-impact-predictor-of-student-success](https://themarkup.org/news/2021/03/02/major-universities-are-using-race-as-a-high-impact-predictor-of-student-success) (2021).
- [11] R. Yu, Q. Li, C. Fischer, S. Doroudi, D. Xu, Towards accurate and fair prediction of college success: evaluating different sources of student data, in: EDM 2020, ERIC, 2020, pp. 292–301.
- [12] F. Marcinkowski, K. Kieslich, C. Starke, M. Lünich, Implications of ai (un-) fairness in higher education admissions: the effects of perceived ai (un-) fairness on exit, voice and organizational reputation, in: ACM FAccT, 2020.
- [13] L. Kondmann, X. X. Zhu, Under the radar—auditing fairness in ml for humanitarian mapping, arXiv preprint arXiv:2108.02137 (2021).
- [14] M. Kearns, S. Neel, A. Roth, Z. S. Wu, Preventing fairness gerrymandering: Auditing and learning for subgroup fairness, in: International Conference on Machine Learning, PMLR, 2018, pp. 2564–2572.
- [15] J. Kleinberg, J. Ludwig, S. Mullainathan, A. Rambachan, Algorithmic fairness, in: Aea papers and proceedings, Vol. 108, 2018, pp. 22–27.
- [16] R. Yu, Q. Li, C. Fischer, S. Doroudi, D. Xu, Towards accurate and fair prediction of college success: evaluating different sources of student data,



- in: Proceedings of the 13th International Conference on Educational Data Mining (EDM 2020), ERIC, 2020, pp. 292–301.
- [17] M. J. Kusner, J. R. Loftus, C. Russell, R. Silva, Counterfactual fairness, arXiv preprint arXiv:1703.06856 (2017).
- [18] G. W. Cole, S. A. Williamson, Avoiding resentment via monotonic fairness, arXiv preprint arXiv:1909.01251 (2019).
- [19] A. Olteanu, C. Castillo, F. Diaz, E. Kiciman, Social data: Biases, methodological pitfalls, and ethical boundaries, *Frontiers in Big Data* 2 (2019) 13.
- [20] S. Barocas, M. Hardt, A. Narayanan, Fairness and machine learning: Limitations and opportunities, [fairmlbook.org](http://fairmlbook.org) (2019).
- [21] A. Asudeh, Z. Jin, H. Jagadish, Assessing and remedying coverage for a given dataset, in: ICDE, 2019, pp. 554–565.
- [22] S. M. Quintana, L. Mahgoub, Ethnic and racial disparities in education: Psychology’s role in understanding and reducing disparities, *Theory Into Practice* 55 (2) (2016).
- [23] N. M. Stephens, S. S. Townsend, A. G. Dittmann, Social-class disparities in higher education and professional workplaces: The role of cultural mismatch, *Current Directions in Psychological Science* 28 (1) (2019) 67–73.
- [24] D. Filmer, The structure of social disparities in education: Gender and wealth, The World Bank, 2000.

- [25] A. Maddrell, K. Strauss, N. J. Thomas, S. Wyse, Mind the gap: Gender disparities still to be addressed in uk higher education geography, *Area* 48 (1) (2016) 48–56.
- [26] D. Voyer, S. D. Voyer, Gender differences in scholastic achievement: A meta-analysis., *Psychological bulletin* 140 (4) (2014) 1174.
- [27] M. Destin, Identity research that engages contextual forces to reduce socioeconomic disparities in education, *Current Directions in Psychological Science* 29 (2) (2020) 161–166.
- [28] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, R. Zemel, Fairness through awareness, in: *ITSC, 2012*, pp. 214–226.
- [29] M. B. Zafar, I. Valera, M. G. Rogriguez, K. P. Gummadi, Fairness constraints: Mechanisms for fair classification, in: *Artificial Intelligence and Statistics*, PMLR, 2017, pp. 962–970.
- [30] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, S. Venkatasubramanian, Certifying and removing disparate impact, in: *SIGKDD, ACM*, 2015, pp. 259–268.
- [31] F. Kamiran, T. Calders, Data preprocessing techniques for classification without discrimination, *Knowledge and Information Systems* 33 (1) (2012) 1–33.
- [32] F. Calmon, D. Wei, B. Vinzamuri, K. N. Ramamurthy, K. R. Varshney, Optimized pre-processing for discrimination prevention, in: *Advances in Neural Information Processing Systems*, 2017, pp. 3992–4001.

- [33] M. B. Zafar, I. Valera, M. G. Rodriguez, K. P. Gummadi, Fairness constraints: Mechanisms for fair classification, arXiv preprint arXiv:1507.05259 (2015).
- [34] H. Zhang, X. Chu, A. Asudeh, S. B. Navathe, Omnifair: A declarative system for model-agnostic group fairness in machine learning, in: SIGMOD, 2021, pp. 2076–2088.
- [35] H. Anahideh, A. Asudeh, S. Thirumuruganathan, Fair active learning, CoRR abs/2006.13025 (2020).
- [36] G. Pleiss, M. Raghavan, F. Wu, J. Kleinberg, K. Q. Weinberger, On fairness and calibration, in: Advances in Neural Information Processing Systems, 2017, pp. 5680–5689.
- [37] M. Zehlike, F. Bonchi, C. Castillo, S. Hajian, M. Megahed, R. Baeza-Yates, Fa\* ir: A fair top-k ranking algorithm, in: CIKM, 2017, pp. 1569–1578.
- [38] I. Žliobaitė, Measuring discrimination in algorithmic decision making, DATA MIN KNOWL DISC 31 (4) (2017) 1060–1089.
- [39] A. Narayanan, Translation tutorial: 21 fairness definitions and their politics, in: ACM FAT\*, 2018.
- [40] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, R. Zemel, Fairness through awareness, in: ITCS, 2012, pp. 214–226.
- [41] M. Hardt, E. Price, N. Srebro, Equality of opportunity in supervised

- learning, in: *Advances in neural information processing systems*, 2016, pp. 3315–3323.
- [42] S. Corbett-Davies, E. Pierson, A. Feller, S. Goel, A. Huq, Algorithmic decision making and the cost of fairness, in: *SIGKDD*, ACM, 2017, pp. 797–806.
- [43] D. Madras, E. Creager, T. Pitassi, R. Zemel, Fairness through causal awareness: Learning causal latent-variable models for biased data, in: *ACM FAT\**, 2019, pp. 349–358.
- [44] K. Makhlouf, S. Zhioua, C. Palamidessi, On the applicability of machine learning fairness notions, *ACM SIGKDD Explorations Newsletter* 23 (1) (2021) 14–23.
- [45] M. Veale, M. Van Kleek, R. Binns, Fairness and accountability design needs for algorithmic support in high-stakes public sector decision-making, in: *ACM CHI*, 2018, pp. 1–14.
- [46] K. Holstein, J. Wortman Vaughan, H. Daumé III, M. Dudik, H. Wallach, Improving fairness in machine learning systems: What do industry practitioners need?, in: *ACM CHI*, 2019, pp. 1–16.
- [47] A. Chouldechova, D. Benavides-Prado, O. Fialko, R. Vaithianathan, A case study of algorithm-assisted decision making in child maltreatment hotline screening decisions, in: *ACM FAT\**, PMLR, 2018, pp. 134–148.
- [48] J. Mokander, L. Floridi, Ethics-based auditing to develop trustworthy ai, *arXiv preprint arXiv:2105.00002* (2021).

- [49] I. D. Raji, A. Smart, R. N. White, M. Mitchell, T. Gebru, B. Hutchinson, J. Smith-Loud, D. Theron, P. Barnes, Closing the ai accountability gap: Defining an end-to-end framework for internal algorithmic auditing, in: ACM FAccT, 2020, pp. 33–44.
- [50] C. Wilson, A. Ghosh, S. Jiang, A. Mislove, L. Baker, J. Szary, K. Trindel, F. Polli, Building and auditing fair algorithms: A case study in candidate screening, in: ACM FAccT, 2021, pp. 666–677.
- [51] J. L. Stephan, E. Davis, J. Lindsay, S. Miller, Who will succeed and who will struggle? predicting early college success with indiana’s student information system. rel 2015-078., Regional Educational Laboratory Midwest (2015).
- [52] M. Ramaswami, R. Bhaskaran, A study on feature selection techniques in educational data mining, arXiv preprint arXiv:0912.3924 (2009).
- [53] T. Chamorro-Premuzic, A. Furnham, Personality, intelligence and approaches to learning as predictors of academic performance, *Personality and individual differences* 44 (7) (2008) 1596–1603.
- [54] O. Zlatkin-Troitschanskaia, J. Schlax, J. Jitomirski, R. Happ, C. Kühling-Thees, S. Brückner, H. Pant, Ethics and fairness in assessing learning outcomes in higher education, *Higher Education Policy* 32 (4) (2019) 537–556.
- [55] D. Avdic, M. Gartell, Working while studying? student aid design and socioeconomic achievement disparities in higher education, *Labour Economics* 33 (2015) 26–40.

- [56] A. Hamoud, A. S. Hashim, W. A. Awadh, Predicting student performance in higher education institutions using decision tree analysis, *International Journal of Interactive Multimedia and Artificial Intelligence* 5 (2018) 26–31.
- [57] K. Pelaez, Latent class analysis and random forest ensemble to identify at-risk students in higher education, Ph.D. thesis, San Diego State University (2018).
- [58] S. A. Dudani, The distance-weighted k-nearest-neighbor rule, *IEEE Transactions on Systems, Man, and Cybernetics* (4) (1976) 325–327.
- [59] T. Tanner, H. Toivonen, Predicting and preventing student failure—using the k-nearest neighbour method to predict student performance in an online course environment, *Int. Journal of Learning Technology* 5 (4) (2010) 356–377.
- [60] E. D. Thompson, B. V. Bowling, R. E. Markle, Predicting student success in a major’s introductory biology course via logistic regression analysis of scientific reasoning ability and mathematics scores, *Research in Science Education* 48 (1) (2018) 151–163.
- [61] M. Agaoglu, Predicting instructor performance using data mining techniques in higher education, *IEEE Access* 4 (2016) 2379–2387.
- [62] P. D. Allison, *Missing data*, Sage publications, 2001.
- [63] D. B. Rubin, Multiple imputation after 18+ years, *Journal of the American statistical Association* 91 (434) (1996) 473–489.

- [64] O. Troyanskaya, M. Cantor, G. Sherlock, P. Brown, T. Hastie, R. Tibshirani, D. Botstein, R. B. Altman, Missing value estimation methods for dna microarrays, *Bioinformatics* 17 (6) (2001) 520–525.
- [65] R. Somasundaram, R. Nedunchezian, Evaluation of three simple imputation methods for enhancing preprocessing of data with missing values, *International Journal of Computer Applications* 21 (10) (2011) 14–19.

## 8. APPENDIX



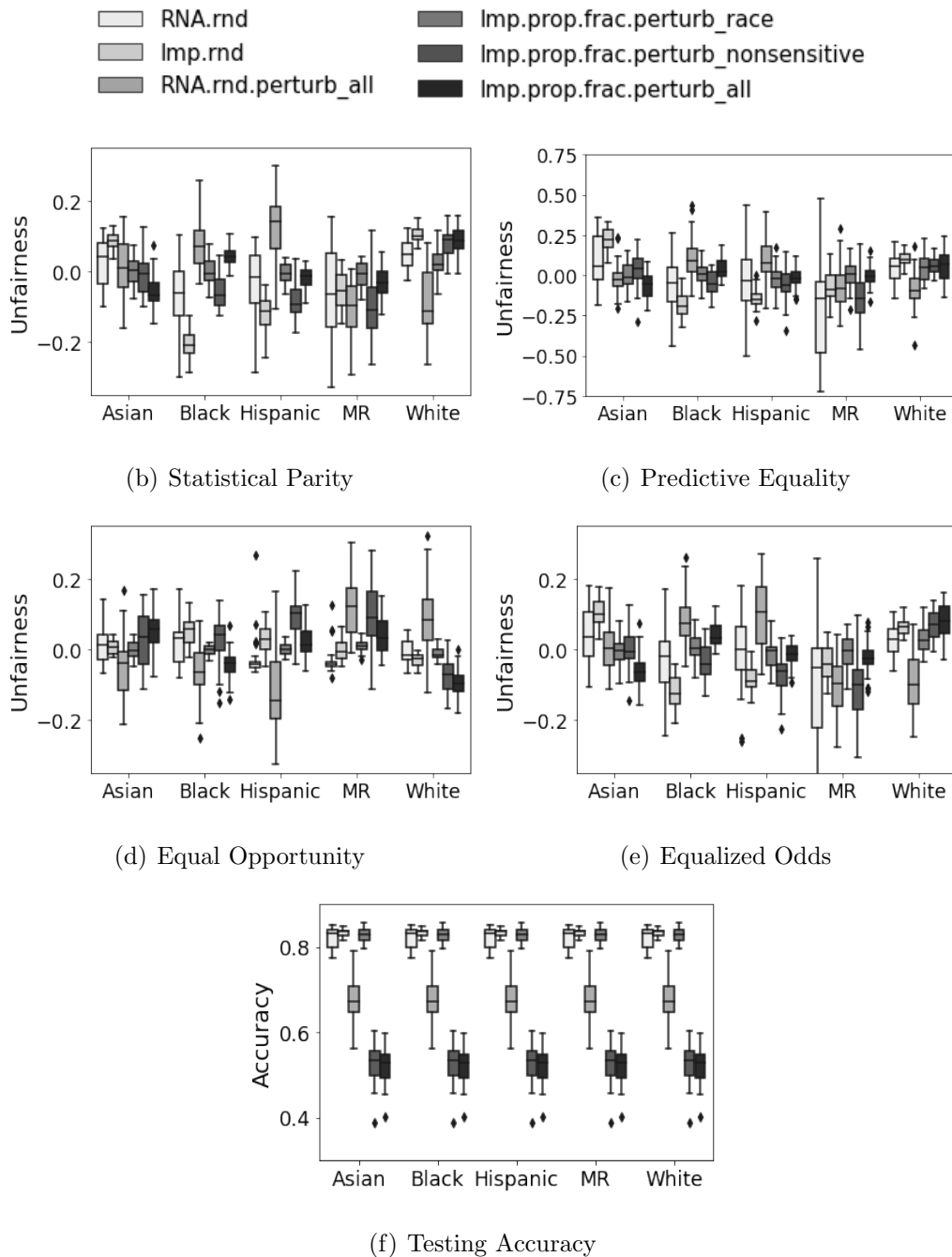


Figure 5: Fairness and Performance Evaluation of Different Perturbation Scenarios compared against Remove-NA and Imp. (SVC)

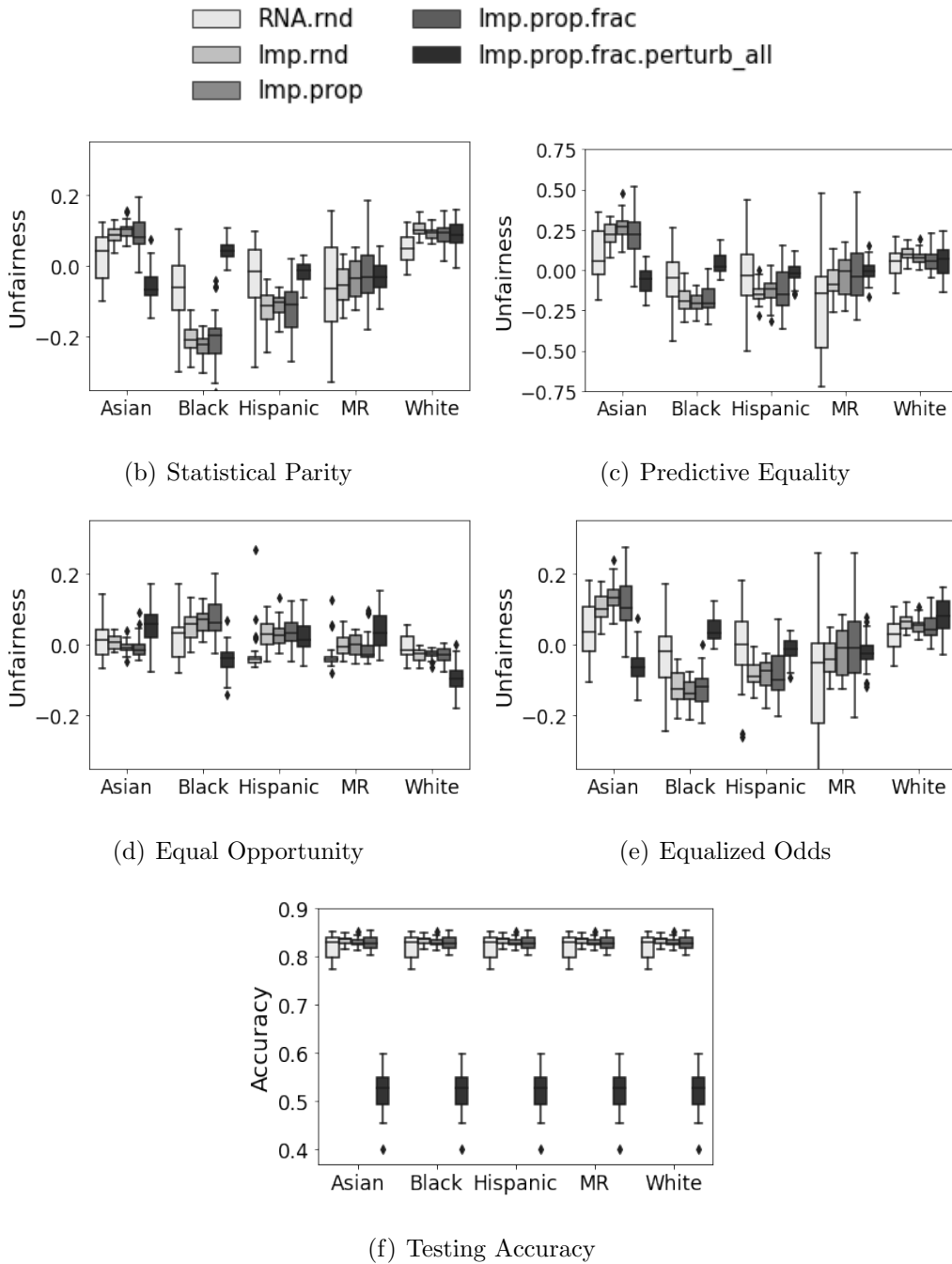


Figure 6: Fairness and Performance Evaluation of Remove-NA, Imputation Scenarios vs. the Perturbation of all attributes (SVC)

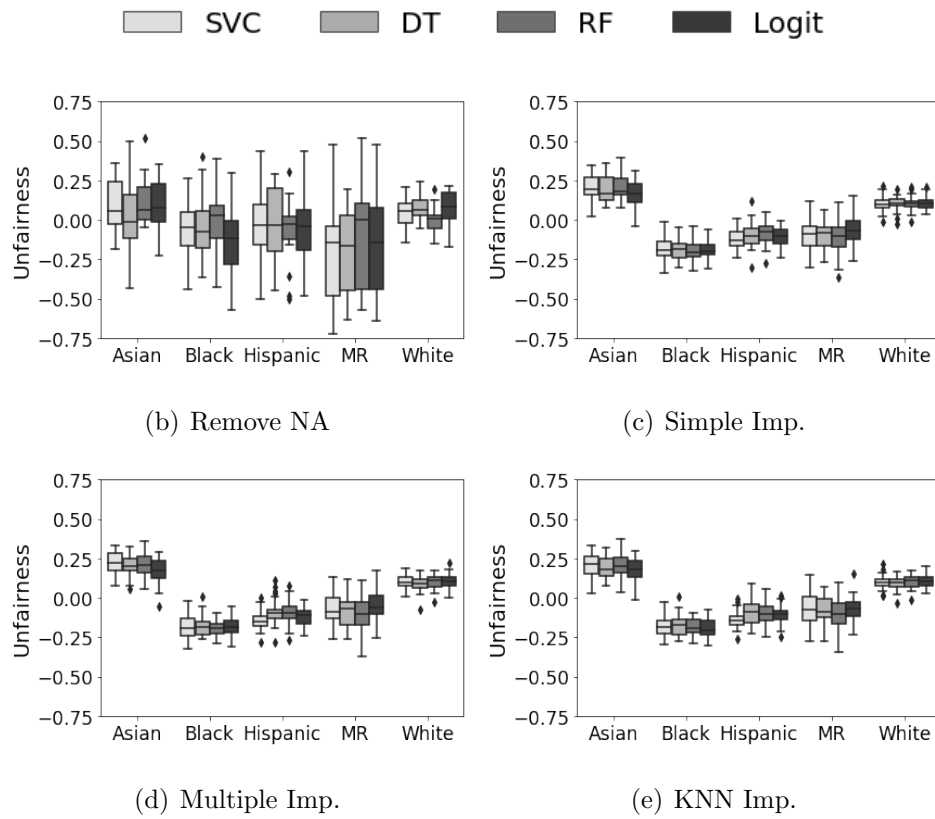


Figure 7: Impact of Imputation - Predictive Equality

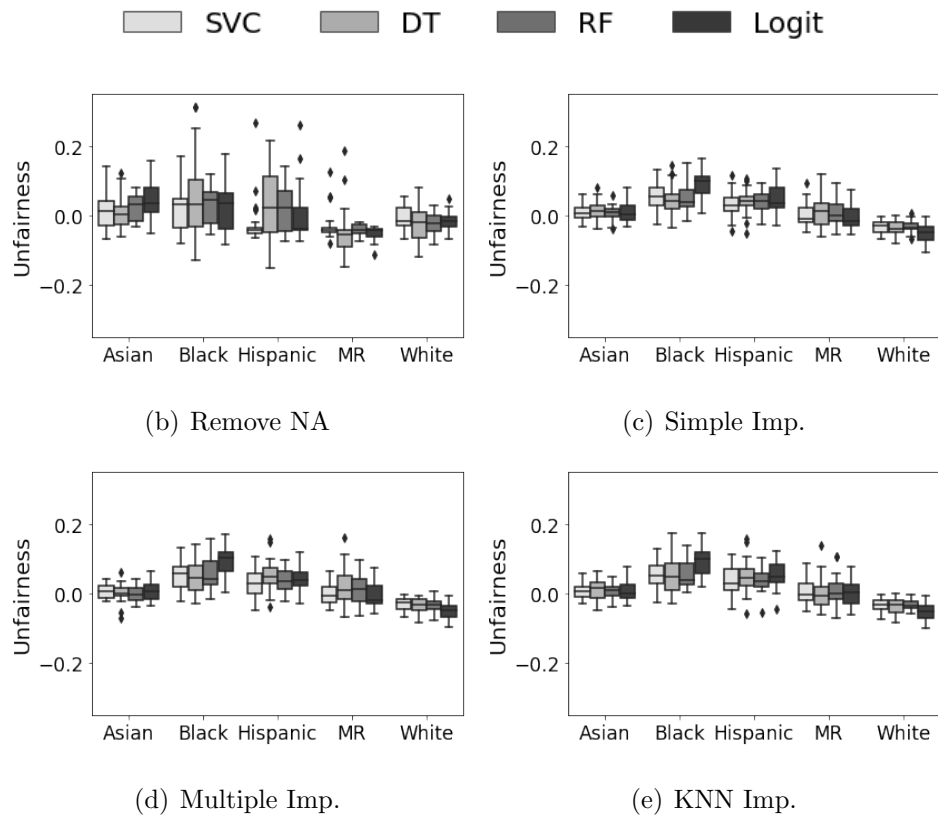


Figure 8: Impact of Imputation- Equal Opportunity

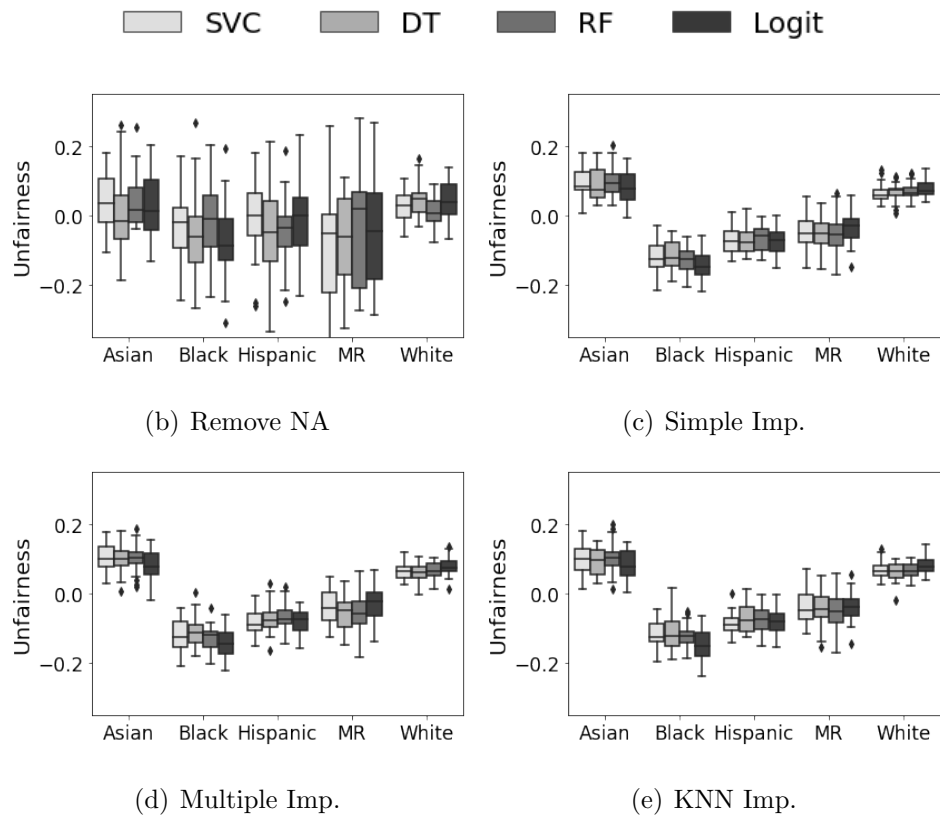


Figure 9: Impact of Imputation- Equalized odds

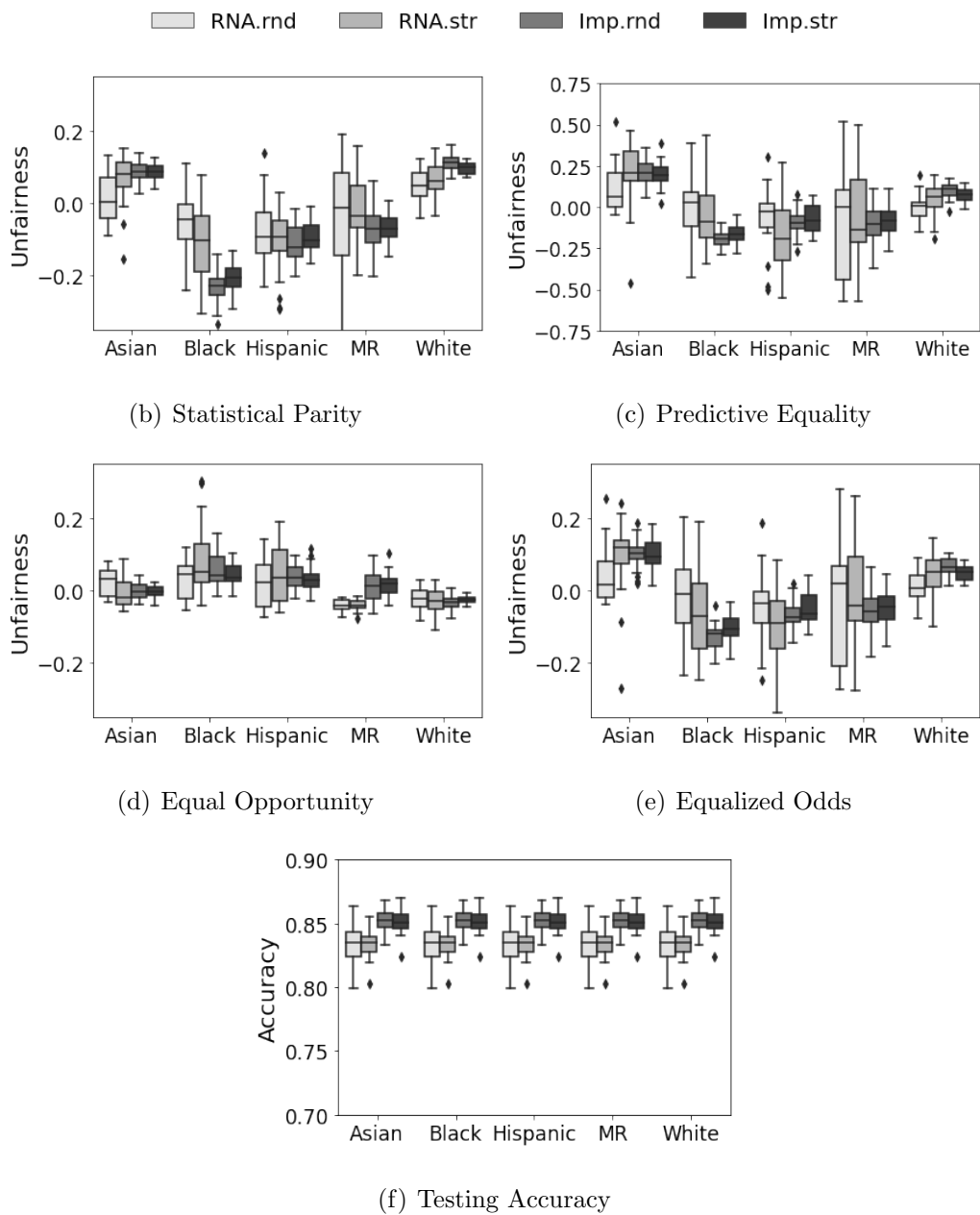


Figure 10: Remove NA vs Imputation (RF)

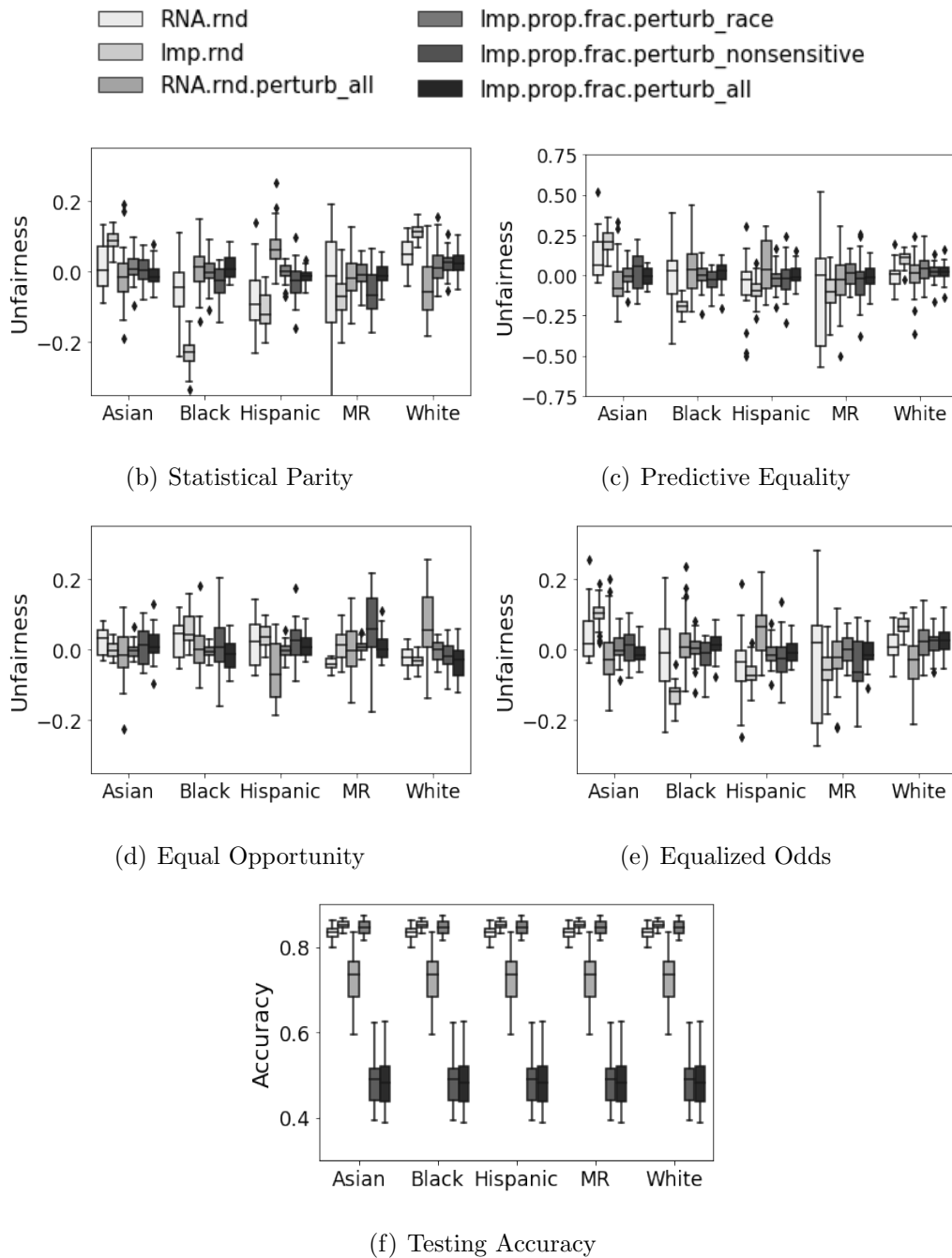


Figure 11: Fairness and Performance Evaluation of Different Perturbation Scenarios compared against Remove-NA and Imp. (RF)

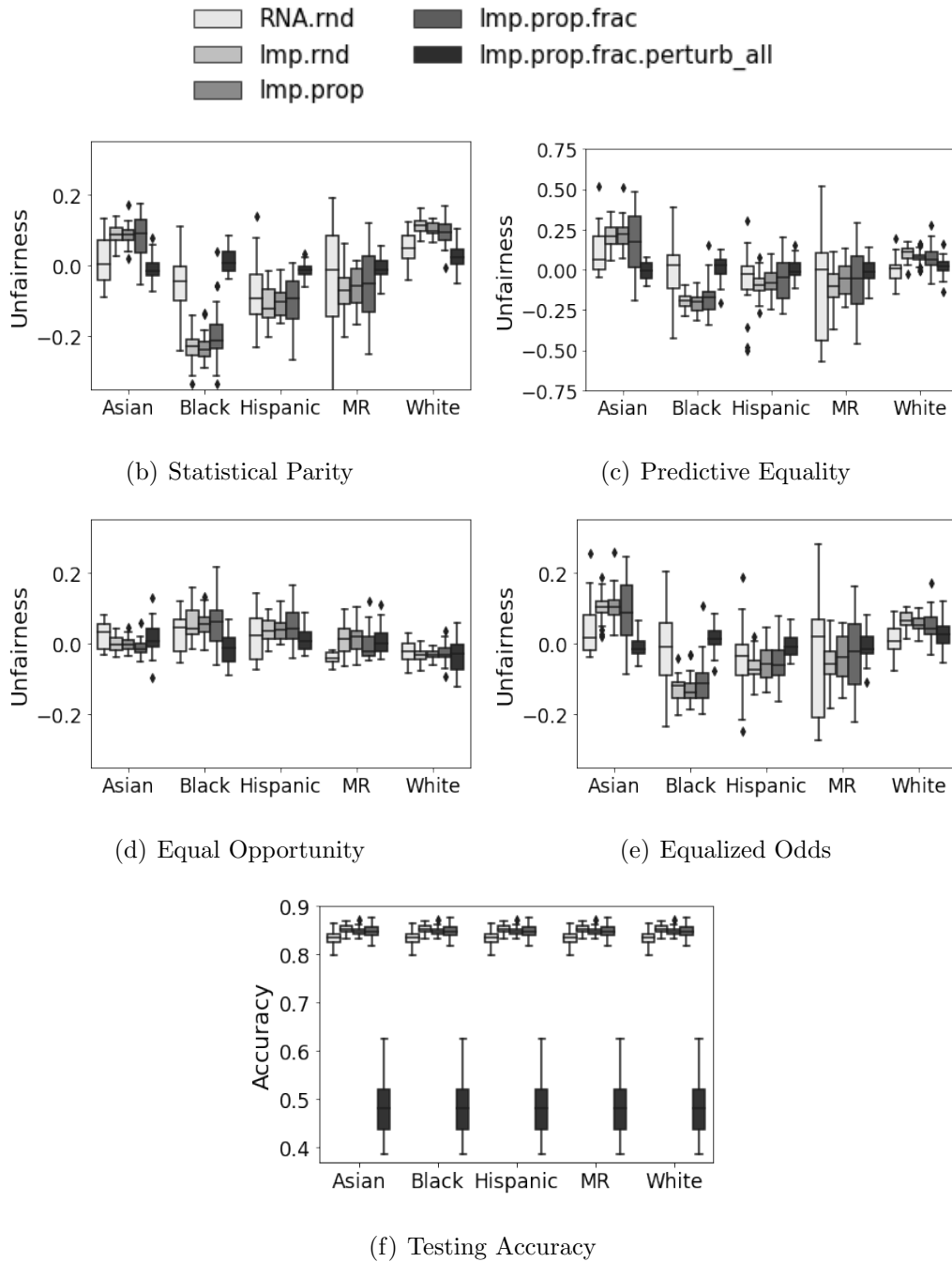


Figure 12: Fairness and Performance Evaluation of Remove-NA, Imputation Scenarios vs. the Perturbation of all attributes (RF)