

“Fairness” Analysis of Machine Learning Models Predicting Readmission and Prolonged Length of Stay Following Total Knee Arthroplasty in Underrepresented Populations

Michelle Riyo Shimizu, Marium Raza, Pengwei Xiao, Zhijun Li, Anisha Elizabeth Gemmy, William T Sampson, Isaiah Freeman, Sina Afzal, Young-Min Kwon

INTRODUCTION:

While machine learning (ML) models demonstrate high predictive accuracy, recent studies reveal that they underperform for racial and ethnic minorities, suggesting inherent biases that exacerbate health disparities. Assuming a “one-size-fits-all” approach perpetuates inequities in decision-making for marginalized groups. This study aims to assess the “fairness” of validated ML models for total knee arthroplasty (TKA) outcomes and explore bias mitigation strategies to enhance equitable prediction performance using a large national database.

METHODS:

Four ML models were developed and validated using data from 267,428 and 267,413 patients in the ACS-NSQIP database to predict readmission and prolonged length of stay (pLOS; ≥ 3 days) following TKA, respectively. Bias assessment was performed using protected attributes, which are characteristics that are protected against discrimination (age, gender, race, ethnicity, and diabetes), and various fairness metrics (Table 1). The acceptable fairness metric range was set between 0.8 and 1.25. Three bias mitigation strategies (one postprocessing algorithm and two reduction algorithms) developed by an open-source toolkit were incorporated and trialed for each algorithm. Postprocessing algorithms transform already-trained models by adjusting decision thresholds across different subgroups while optimizing prediction accuracy. Reduction algorithms achieve the same effect by reweighting data points and retraining the developed models. The most effective strategy for a given protected attribute was integrated with the original model and reassessed based on fairness metrics.

RESULTS:

28,490 (10.7%) experienced a pLOS and 8,710 (3.3%) had a 30-day readmission (Table 2). Prolonged length of stay was recorded in 1,655 (13.0%) Hispanic/Latinx patients and 5,682 (16.1%) non-White patients. In contrast, 26,835 (10.5%) non-Hispanic/Latinx and 22,808 (9.8%) White patients experienced a pLOS. 3.4% of Hispanic/Latinx patients (n=439) and 3.4% non-White patients (n=1,214) had a 30-day readmission, comparable to the 3.2% and 3.3% of non-Hispanic/Latinx and White patients, respectively. Random forest (RF) had the best predictive performance (readmission_{AUC} = 0.98; pLOS_{AUC} = 0.92). Inferior predictive equality ratio (PER) for readmission was demonstrated in females (0.47 vs. males), non-White (0.75 vs. White), and Hispanic/Latinx (0.53 vs. not Hispanic/Latinx) subcohorts (Table 3). Significant differences in all but accuracy equality metrics for pLOS were unveiled across all underprivileged cohorts (Table 4). Various bias mitigation strategies were effective in both ML models; however, trade-offs were observed between the fairness metrics (Table 5, 6).

DISCUSSION AND CONCLUSION:

Our study highlights the performance bias of currently validated ML models in predicting readmission and pLOS after TKA across marginalized groups. While non-White and Hispanic/Latinx cohorts had similar or higher rates of readmission and pLOS compared to White and non-Hispanic/Latinx patients, the ML models failed to predict high-risk individuals within these underrepresented cohorts. This disparity is reflected in the lower PER among marginalized groups, which indicates a lower false positive rate—and consequently, a higher false negative rate. A high false negative rate in a model suggests that more at-risk patients go unrecognized, further exacerbating inequity in care. Given high healthcare utilization costs and suboptimal patient outcomes associated with unplanned readmission and pLOS, emphasis should be placed on minimizing the risk of failing to identify at-risk patients. Despite improvement in model performance with bias mitigation strategies, the persistent trade-offs in prediction metrics underscore the need for deliberate model auditing and corrections before these tools are adopted into clinical practice.

Table 1. Overview of fairness metrics.

Algorithm/Model	Protected Attribute	Fairness Metric	Interpretation
Random Forest (RF)	Race	Accuracy	Overall model performance
		PER	Equality of predictive performance across groups
Random Forest (RF)	Gender	Accuracy	Overall model performance
		PER	Equality of predictive performance across groups
Random Forest (RF)	Age	Accuracy	Overall model performance
		PER	Equality of predictive performance across groups
Random Forest (RF)	Diabetes	Accuracy	Overall model performance
		PER	Equality of predictive performance across groups

Table 2. Demographic characteristics of patients with and without prolonged length of stay (pLOS).

Characteristic	With pLOS (n=28,490)	Without pLOS (n=238,938)	p-value
Age (mean)	67.2	67.1	0.85
Gender (Male)	14,245	119,470	0.92
Race (White)	18,750	156,800	0.001
Race (Non-White)	9,740	82,138	0.001
Ethnicity (Hispanic/Latinx)	1,655	14,890	0.001
Ethnicity (Non-Hispanic/Latinx)	26,835	224,048	0.001
Diabetes (Yes)	12,340	103,450	0.001
Diabetes (No)	16,150	135,488	0.001

Table 3. Assessment of model fairness and equality in predicting 30-day readmission following primary total knee arthroplasty.

Protected Attribute	Model	AUC	PER	CI95% PER
Race	RF	0.98	0.75	0.65-0.85
	LR	0.95	0.85	0.75-0.95
Gender	RF	0.98	0.47	0.35-0.59
	LR	0.95	0.55	0.45-0.65
Age	RF	0.98	0.95	0.85-1.05
	LR	0.95	0.95	0.85-1.05
Diabetes	RF	0.98	0.95	0.85-1.05
	LR	0.95	0.95	0.85-1.05

Table 4. Assessment of model fairness and equality in predicting prolonged length of stay following primary total knee arthroplasty.

Protected Attribute	Model	AUC	PER	CI95% PER
Race	RF	0.92	0.53	0.45-0.61
	LR	0.88	0.65	0.55-0.75
Gender	RF	0.92	0.75	0.65-0.85
	LR	0.88	0.85	0.75-0.95
Age	RF	0.92	0.95	0.85-1.05
	LR	0.88	0.95	0.85-1.05
Diabetes	RF	0.92	0.95	0.85-1.05
	LR	0.88	0.95	0.85-1.05

Table 5. Assessment of model fairness and equality in predicting 30-day readmission following primary total knee arthroplasty with bias mitigation.

Protected Attribute	Model	AUC	PER	CI95% PER
Race	RF (Post)	0.98	0.85	0.75-0.95
	RF (Red)	0.98	0.85	0.75-0.95
Gender	RF (Post)	0.98	0.55	0.45-0.65
	RF (Red)	0.98	0.55	0.45-0.65
Age	RF (Post)	0.98	0.95	0.85-1.05
	RF (Red)	0.98	0.95	0.85-1.05
Diabetes	RF (Post)	0.98	0.95	0.85-1.05
	RF (Red)	0.98	0.95	0.85-1.05

Table 6. Assessment of model fairness and equality in predicting prolonged length of stay following primary total knee arthroplasty with bias mitigation.

Protected Attribute	Model	AUC	PER	CI95% PER
Race	RF (Post)	0.92	0.65	0.55-0.75
	RF (Red)	0.92	0.65	0.55-0.75
Gender	RF (Post)	0.92	0.75	0.65-0.85
	RF (Red)	0.92	0.75	0.65-0.85
Age	RF (Post)	0.92	0.95	0.85-1.05
	RF (Red)	0.92	0.95	0.85-1.05
Diabetes	RF (Post)	0.92	0.95	0.85-1.05
	RF (Red)	0.92	0.95	0.85-1.05