## Machine Learning Can Mitigate Racial Bias in Orthopaedics: Predicting Overnight Stay After Knee Arthroscopy

Jonathan S. Lee<sup>1</sup>, Bilal S. Siddiq<sup>1</sup>, Stephen M. Gillinov<sup>1</sup>, Kieran S. Dowley<sup>1</sup>, Megan J. Wei<sup>2</sup>, Nathan J. Cherian<sup>1</sup>, Mark P. Cote<sup>1</sup>, Scott D. Martin<sup>1</sup>

<sup>1</sup>Massachusetts General Hospital Division of Sports Medicine, Boston, MA, <sup>2</sup>Johns Hopkins Department of Computer Science, Baltimore, MD

Email of Presenting Author: <u>jlee376@mgh.harvard.edu</u>

**INTRODUCTION**: Over the past 3 years, the utilization of machine learning in orthopaedics has risen exponentially. Past studies have employed machine learning to predict hospital readmission, rates of achieving clinical thresholds, and functional improvements. Due to the utility of these models, some orthopaedic researchers have gone as far as publishing these algorithms through online calculators, allowing users to input pertinent patient data and derive a predicted clinical outcome. Despite machine learning's immense potential for clinical application, before researchers can offer these tools to a broad and diverse patient population, they must first ensure that their machine learning model performs with high degrees of accuracy irrespective of patient race. The purpose of the present study was to develop a racially equitable machine learning model that accurately predicts overnight stay following arthroscopic knee surgery.

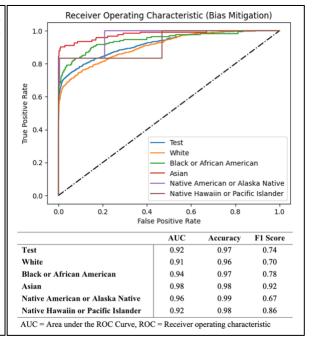
METHODS: This retrospective study queried the NSQIP database for patients ≥ 18 years old who underwent knee arthroscopy between 2011 and 2020. Patients who self-identified as White, Black or African American, Asian, Native American or Alaskan Native, or Native Hawaiian or Pacific Islander and had complete demographic, intraoperative, and post-operative outcome data were included in our analysis. All model development and statistical analysis was carried out in Python using packages from Scikit-learn. Predictive analysis was performed using a Random Forest (RF) and performance was optimized through feature selection, hyperparameter tuning with cross validation, and over- and under-sampling techniques with SMOTE. The following patient variables to train the RF model included but were not limited to sex, age, BMI, wound classification, ASA class, and type of anesthesia. Patients with a hospital length of stay (LOS) ≥ 1 day were placed into the overnight stay cohort, while patients with a LOS < 1 represented the same-day discharge cohort. Model performance was assessed using AUC, accuracy, and F1 score.

**RESULTS:** Of the 75,775 patients who met inclusion criteria, 7.5% (n=5,714) required overnight stay. When stratified by racial self-identification, most patients identified as White (78.3%), followed by Black or African American (13.7%), Asian (5.4%), Native American or Native Alaskan (1.5%), and Native Hawaiian or Pacific Islander (1.1%) (**Table 1**). Despite a lack of representation among minority patient groups, the racial-equitably trained RF model showed high levels of performance for AUC (mean =  $0.94 \pm 0.03$ ) and accuracy (mean =  $0.97 \pm 0.01$ ) irrespective of patient race. F1 score, however, varied with the model performing best for Asian (0.92), Native Hawaiian or Pacific Islander (0.86), and Black or African American patients (0.78) (**Figure 1**). Of the variables included to train the RF model, 5 features accounted for 73% of model development (**Table 2**).

**DISCUSSION**: We developed a racially equitable machine learning model that effectively predicted overnight stay in patients undergoing knee arthroscopy. As the utilization of machine learning in orthopaedics continues to rise, it is crucial that models intended for clinical application are generalizable and exhibit high levels of performance for minority patients.

SIGNIFICANCE/CLINICAL RELEVANCE: The onus is on orthopaedic researchers to assess the internal validity of their machine learning models before offering them to a broader patient population. Failing to do so may result in inequitable clinical models that exacerbate poor health outcomes among minority patients.

	Patients (n=75,775)
Same-Day Discharge	5,714 (7.5%)
Overnight Stay	70,061 (92.5%)
Age, years	$37.2 \pm 14.4$
BMI, kg/m <sup>2</sup>	$29.5 \pm 7.3$
Sex	
Male	32,101 (42.4%
Female	43,674 (57.6%
Racial Self-Identification	
White	59,313 (78.3%
Black or African American	10,359 (13.7%
Asian	4,099 (5.4%
Native American or Alaska Native	1,152 (1.5%
Native Hawaiin or Pacific Islander	852 (1.1%
Current Procedural Terminology	
29875	5,291 (7.0%
29876	6,819 (9.0%
29877	15,361 (20.3%
29882	9,025 (11.9%
29888	39,279 (51.8%



Variables	Contribution
Arthroscopic knee procedure	36%
Admission quartile	16%
Anesthesia type	13%
ASA class	5%
Days from hospital admission to operation	3%