

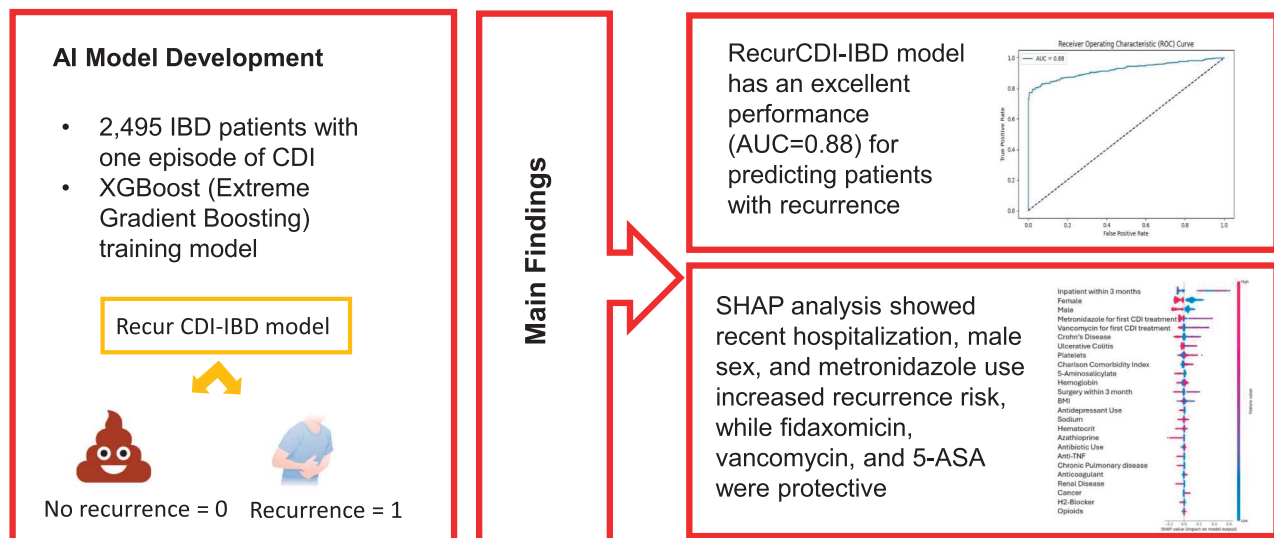
# Machine Learning Model Predicts Recurrent *Clostridioides difficile* Infection in Patients With Inflammatory Bowel Disease (Recur CDI-IBD)

Raseen Tariq, MD, MS<sup>1,2</sup>, Ankita Sethi, BS<sup>1</sup>, Shivaram Arunachalam, PhD<sup>1</sup>, Darrell S. Pardi, MD, MS, FACP<sup>1</sup>, William A. Faubion, MD<sup>3</sup> and Sahil Khanna, MBBS, MS, FACP<sup>1</sup>

**INTRODUCTION:** *Clostridioides difficile* infection (CDI) presents a significant challenge in patients with inflammatory bowel disease (IBD), with high recurrence rates and complications. Predicting recurrent CDI (rCDI) in patients with IBD is crucial for implementing targeted interventions to improve patient outcomes. This study aimed to develop and validate a predictive model (RecurCDI-IBD) using supervised machine learning to identify patients with IBD at high risk of developing rCDI.

**METHODS:** Data were collected from adult patients with IBD diagnosed with CDI between 2013 and 2021. Inclusion criteria included adult patients with a confirmed diagnosis of CDI and a history of IBD. The Gradient Boosting Machine learning model (XGBoost) was used to train a binary classification model. Feature engineering included demographic data (age and sex), clinical data (IBD subtype, medication use, and comorbidities), and laboratory data. The primary outcome was the occurrence of rCDI within 60 days of the initial CDI episode.

## Machine Learning Model Predicts Recurrent *Clostridioides difficile* Infection in Patients With Inflammatory Bowel Disease (Recur CDI-IBD)



Tariq R et al. *Am J Gastroenterol*. 2026. doi: 10.14309/ajg.0000000000003922  
 © 2026 by The American College of Gastroenterology

**AJG** The American Journal of GASTROENTEROLOGY

<sup>1</sup>Division of Gastroenterology and Hepatology, Mayo Clinic, Rochester, Minnesota, USA; <sup>2</sup>Division of Gastroenterology and Hepatology, Virginia Commonwealth University, Richmond, Virginia, USA; <sup>3</sup>Division of Gastroenterology and Hepatology, Mayo Clinic, Arizona, Phoenix, USA. **Correspondence:** Sahil Khanna, MBBS, MS, FACP. E-mail: khanna.sahil@mayo.edu.

Received July 2, 2025; accepted January 15, 2026; published online January 21, 2026

**RESULTS:** The RecurCDI-IBD model achieved an accuracy of 80.05% and an Area Under the Curve of 0.88 for predicting rCDI. Key predictive features included IBD subtype, sex, specific medications (such as steroids and anti-TNF agents), and comorbidities (such as chronic pulmonary and renal disease).

**DISCUSSION:** The RecurCDI-IBD model demonstrates good discriminatory ability with balanced precision and recall in identifying patients with IBD at higher risk for rCDI. These findings highlight the potential of data-driven approaches to support clinical risk assessment. Further studies incorporating larger and more diverse cohorts and prospective external validation are needed to confirm generalizability and optimize clinical applicability.

**KEYWORDS:** *C. difficile*; machine learning; inflammatory bowel disease; prediction

**ABBREVIATIONS:** 5-ASA, 5-aminosalicylic acid; AUC, Area under the receiver operating characteristic curve; BMI, Body mass index; CD, Crohn's disease; CDI, Clostridioides difficile infection; CRP, C-reactive protein; EHR, Electronic health record; F1-score, Harmonic mean of precision and recall; IBD, Inflammatory bowel disease; ICD, International Classification of Diseases; PCR, Polymerase chain reaction; rCDI, Recurrent Clostridioides difficile infection; SHAP, SHapley Additive exPlanations; SMOTE, Synthetic Minority Over-sampling Technique; UC, Ulcerative colitis; XGBoost, Extreme Gradient Boosting

**SUPPLEMENTARY MATERIAL** accompanies this paper at <http://links.lww.com/AJG/D866>, <http://links.lww.com/AJG/D867>, <http://links.lww.com/AJG/D868>

*Am J Gastroenterol* 2026;121:942–949. <https://doi.org/10.14309/ajg.0000000000003922>

## INTRODUCTION

*Clostridioides difficile* infection (CDI) remains a major challenge in patients with inflammatory bowel disease (IBD), frequently precipitating disease flares, treatment failure, hospitalizations, and surgery (1). Recurrence after initial episode is common and represents a major unmet need, leading to repeated antibiotics exposure, prolonged morbidity, and impaired quality of life. Predicting which patients with IBD are most likely to experience rCDI is therefore critical to guide timely, preventive care and reduce downstream complications.

Despite the clinical significance of rCDI, current prediction strategies remain limited. The ability to predict rCDI in patients with IBD remains inadequate because of complexity of both diseases, variability in patient populations, and the multifaceted nature of both CDI and IBD (2,3). The complexity arises from differences in immune response, gut microbiota composition, genetic predispositions, and external factors such as diet and antibiotic use. In addition, the effectiveness of treatments can vary, and patient variability in demographics, comorbidities, and adherence to treatment regimens further complicates prediction efforts (4,5). As a result, no validated predictive tool currently exists to accurately stratify rCDI risk in IBD, highlighting a key gap in existing research.

Supervised machine learning (ML) has emerged as a useful research tool, specifically where traditional statistical methods fall short (6). The complexity of CDI and IBD results in nonlinear relationships and high-dimensional data that traditional statistical models struggle to handle effectively. Machine learning models can manage these complexities by identifying patterns and interactions that are not apparent through conventional methods (7). They can integrate diverse data types and continuously learn from new data, making them well-suited for predicting rCDI in patients with IBD (8,9).

The aim of this study was to develop and internally validate a supervised machine learning model (Recur CDI-IBD) to predict rCDI among patients with IBD.

## METHODS

### Population

We used data from adult patients (aged 18 or older) with a confirmed diagnosis of IBD who developed CDI between 2013 and

2021. IBD was defined using the *International Classification of Diseases (ICD-9 and ICD-10)* codes, requiring at least 2 IBD-related codes recorded at least 30 days apart (see Appendix 1, Supplementary Digital Content, <http://links.lww.com/AJG/D866>). One of these codes was required to be from an outpatient setting, a method validated in previous studies with a positive predictive value of 0.83 for Crohn's disease (CD) and 0.89 for ulcerative colitis (UC) (10). To further characterize the IBD cohort, we extracted additional variables capturing treatment exposure, disease subtype, and surgical history. IBD medications included corticosteroids, immunomodulators (azathioprine and methotrexate), biologics (anti-tumor necrosis factor, vedolizumab, and ustekinumab), small-molecule therapy (tofacitinib), and 5-aminosalicylates (5-ASA). Additional gastrointestinal medications including bile acid binders, antidiarrheals, and antimotility agents were included to capture supportive or adjunctive therapies. Surrogates of disease activity were incorporated using recent corticosteroid exposure, hospitalization within 3 months before CDI, and laboratory markers including hemoglobin, hematocrit, and platelet count. C reactive protein and fecal calprotectin values were missing in more than 40% of patients and were therefore excluded from the final analysis. In addition, a surgical history variable identified patients with prior bowel resection, while IBD subtypes (CD vs. UC) were coded as separate variables (see Appendix 2, Supplementary Digital Content, <http://links.lww.com/AJG/D867>).

### CDI diagnosis criteria

Incident CDI was identified using ICD codes, a positive stool test for *C. difficile* detected by polymerase chain reaction, and prescription of vancomycin, metronidazole, or fidaxomicin within 3 days of the positive test. To ensure these were truly new episodes, we excluded any patients who had a positive CDI test or received CDI treatment in the preceding 60 days. This 60-day exclusion window reflects standard clinical definitions of CDI recurrence.

### rCDI definition (target variable)

Our outcome of interest was rCDI, defined as either a positive repeat polymerase chain reaction test with initiation of a new

course of CDI-directed antibiotics (vancomycin, fidaxomicin, or metronidazole) or initiation of a new CDI-directed antibiotic course alone within 60 days of completing therapy for the index episode (11). Because some patients with suspected recurrence are treated empirically without repeat testing, antibiotic initiation alone was also considered indicative of rCDI to reflect real-world clinical practice and maximize the model's applicability.

### Analysis methods

We applied a supervised ML approach using the XGBoost (Extreme Gradient Boosting) algorithm (12) to train a predictive model for rCDI (rCDI). This model was designed to reflect real-world clinical data inputs and to support clinicians in identifying patients with IBD at highest risk of recurrence. Features (columns) with more than 40% missing data were excluded to ensure clinical reliability. For the remaining missing values, we used zero imputation—replacing missing categorical values with zeros—a pragmatic method aligned with standard data handling in clinical informatics (13).

Given the imbalance in recurrence outcomes, we applied the synthetic minority oversampling technique to balance the data set. The synthetic minority oversampling technique creates synthetic examples in the minority class, helping the model detect high-risk patients without overfitting to a small subgroup (14). This step is crucial for maintaining clinical utility, particularly when rare but impactful outcomes such as rCDI are the target. Additional details on model selection, missing data handling, and class-balancing techniques are provided in Appendix 3 (see the Supplementary Digital Content, <http://links.lww.com/AJG/D868>).

### Feature engineering and data set

We extracted features commonly available in the electronic health record (EHR), making the model translatable to clinical practice. These included demographics (age, sex, race, and BMI), IBD subtype, Charlson Comorbidity Index, recent hospitalization, medication exposures (corticosteroids, biologics, immunomodulators, 5-ASAs, and antibiotics), and laboratory values within 3 months of CDI. Medication features were binary (yes/no exposure), while laboratory values were treated as continuous variables. We used one-hot encoding for categorical data to accommodate the input format of the ML model. Of 141 candidate features, 42 were removed because of missingness, leaving 99

features for training (see Appendix 2, Supplementary Digital Content, <http://links.lww.com/AJG/D867>).

### Model training and evaluation

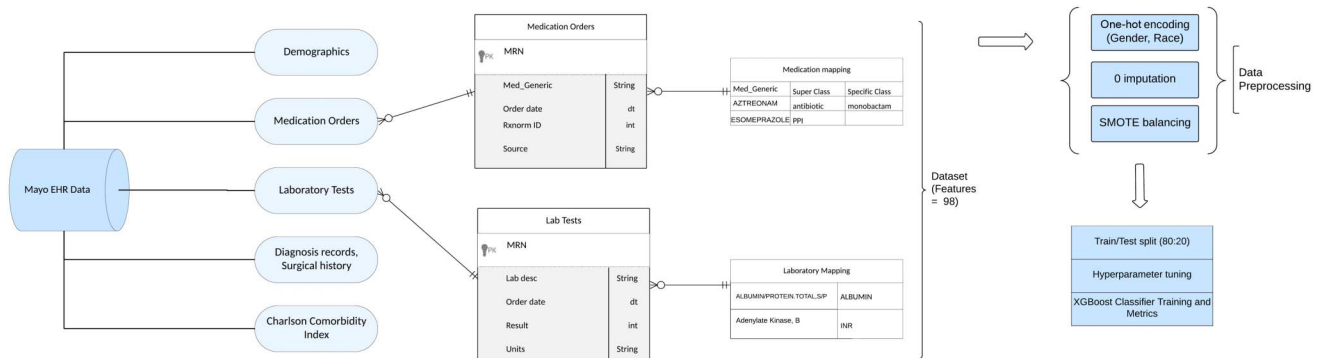
We randomly split the cohort into training (80%) and testing (20%) sets. The XGBoost algorithm was trained on the former, with hyperparameter tuning performed using grid search and 5-fold cross-validation. This method ensured the model was generalizable and not overfitted to any particular subset of patients (Figure 1).

To evaluate the performance of our model, we applied several clinically relevant metrics. Accuracy represented the overall proportion of correct predictions and served as a basic indicator of model performance. However, because rCDI is a relatively infrequent outcome, we also assessed the area under the receiver operating characteristic curve (AUC), which reflects the model's ability to discriminate between patients who will and will not experience recurrence across a range of thresholds (15). Sensitivity, or recall, measured the proportion of actual recurrent cases the model correctly identified—crucial for minimizing missed diagnoses in high-risk patients. Precision captured how often the model's predictions of recurrence were accurate, helping to avoid unnecessary treatment. Finally, the F1-score, a harmonic mean of sensitivity and precision, provided a balanced assessment especially suited for imbalanced data sets, where 1 outcome (non-recurrence) is more common.

Performance was evaluated using both binary classification and predicted probabilities. AUC in particular was emphasized to account for clinical scenarios where a risk estimate (rather than binary classification) may guide patient management.

### Model interpretability

To support clinician adoption, we prioritized transparency. We assessed feature contributions using XGBoost's internal feature importance and SHapley Additive Explanations (SHAP). XGBoost feature (16) importance identified which variables were most influential in decision-making across all patients. SHAP provided individualized risk explanations by measuring the marginal contribution of each feature for every prediction (17). These methods enable gastroenterologists to not only see which patients are high-risk, but understand why—thereby enhancing trust, clinical reasoning, and shared decision-making (Figure 1).



**Figure 1.** Flowsheet for machine learning model development for predicting recurrent CDI. EHR, electronic health record; INR, international normalized ratio; MRN, medical record number; PPI, proton pump inhibitors; SMOTE, synthetic minority oversampling technique.

**Table 1.** Baseline characteristics of patients included in model

Characteristic	Total patients	No recurrence	Recurrence
Total patients	2,495	1,841	654
Sex (female)	1,376	1,015	361
White	2,302	1,698	610
Median age (yr)	50.4 (18–87)	48.5 (18–87)	53.9 (18–86)
Median BMI	25.5 (11.4–56)	25.7 (11.7–55.2)	25.1 (11.7–55)
Inpatient at the time of first CDI	739	489	250
Medications			
5-ASA	535	422	113
Acetaminophen	1,231	839	392
Antacid	133	84	49
Antibiotic	1,596	1,102	494
Anticoagulant	901	595	306
Antidepressant	647	462	185
Anti-TNF	310	245	65
Ulcerative colitis	1,797	1,306	491
Crohn's disease	647	503	144
Comorbidities			
Myocardial infarction	49	31	18
CHF	130	82	48
Peripheral vascular disease	141	87	54
Cerebrovascular disease	74	50	24
Dementia	29	18	11
Chronic pulmonary disease	270	186	84
Diabetes	205	124	81
Other characteristics			
Undergone surgery in 3 mo up to the incident CDI	219	166	53
Charlson Comorbidity Index	1 (0–19)	1 (0–19)	2 (0–16)
Hematocrit within 3 mo	37.6 (17.2–49.1)	38 (17.2–49)	36.2 (18.7–48)
Hemoglobin	12.2 (4.8–18.1)	12.4 (4.8–18.1)	11.8 (5.8–16.8)
Platelets	267 (6–583)	266 (9–516)	272 (6–583)

5-ASA, 5-aminosalicylates; BMI, body mass index; CDI, *Clostridioides difficile* infection; CHF, congestive heart failure; TNF, tumor necrosis factor.

## RESULTS

### Patient characteristics

The study included a total of 2,495 patients of whom 1,841 had no recurrence of CDI and 654 experienced recurrences. Among them, 55.1% were female, the median age of the patients was 50.4 years, ranging from 18 to 87 years. In the overall cohort of 2,495 patients, 72.0% had UC and 25.9% had CD. At the time of the first CDI, 739 patients were inpatients, with 489 in the no recurrence group and 250 in the recurrence group. Table 1 summarizes the details of characteristics of patients included in the model.

### Model performance

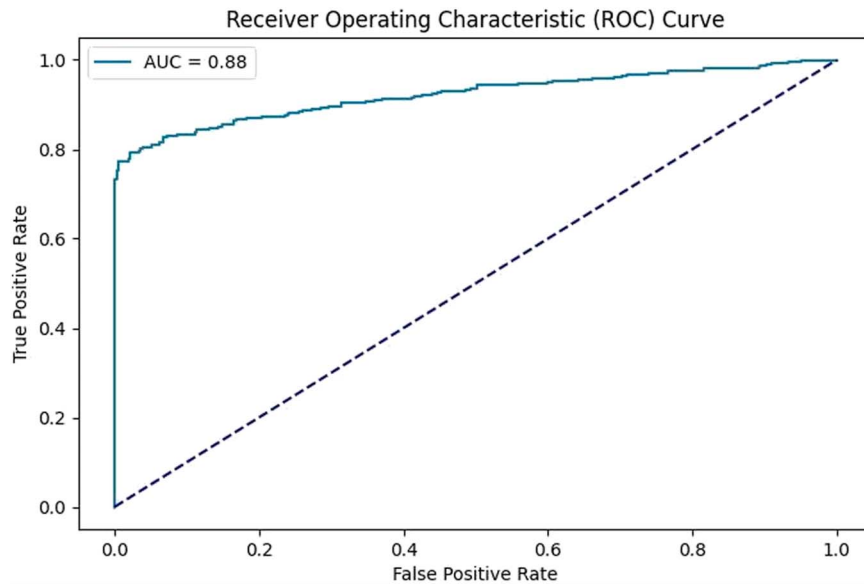
The performance evaluation of the Recur CDI-IBD model, based on independent test set, yielded an accuracy of 80.05% from predicted labels, with an AUC of 0.88 based on predicted probabilities,

indicating a strong ability to discriminate between recurrent and non-rCDI cases. The model demonstrated a sensitivity (recall) of 0.76, a specificity of 0.84, and a precision (positive predictive value) of 0.83 and negative predictive value of 0.77 (Figure 2).

The detailed classification report showed a precision of 0.78, recall of 0.84, and F1-score of 0.81 for class 0 (no recurrence), and a precision of 0.83, recall of 0.76, and F1-score of 0.79 for class 1 (recurrence).

### Mean performance metrics across 5-fold cross validation

The Recur CDI-IBD model demonstrated consistent performance in predicting CDI recurrence within 60 days, as reflected in the mean performance metrics across multiple runs (Table 2). These results are derived from 5-fold cross-validation performed on the training data in which the data set was randomly partitioned into 5 equal subsets;



**Figure 2.** ROC curve for final optimized model evaluated on the independent test set (AUC = 0.88). AUC, area under the receiver operating characteristic curve.

in each iteration, 4 folds were used for training and 1 for validation, and model performance was averaged across all folds. The model achieved a mean accuracy of 72.75% and an AUC of 0.82, indicating moderate discriminatory power. The model's balanced sensitivity (0.78) and specificity (0.75) suggests it can distinguish recurrent from nonrecurrent cases. With a precision of 0.72 and recall of 0.71, the Recur CDI-IBD model demonstrates promising performance for identifying CDI recurrence within 60 days and warrant validation in external cohorts (Figure 3).

#### XGBoost feature importance

The feature importance plot from the XGBoost model highlights key predictors for assessing the risk of rCDI in patients with IBD. The top feature, UC, indicates a significant risk factor for CDI recurrence. Sex-specific impacts were notable, with male sex showing a higher importance than female sex. The use of 5-aminosalicylic acid significantly influenced recurrence risk. The top 25 features are shown in Figure 4.

#### SHAP analysis of top features

The SHAP summary plot (Figure 5) illustrates the magnitude and direction of the top predictive features influencing rCDI risk in patients with IBD. Recent inpatient status within 3 months was the most impactful predictor, with higher SHAP values indicating a strong positive association with recurrence risk, likely reflecting more severe baseline illness or healthcare exposure. Male sex and initial metronidazole therapy similarly shifted model predictions toward recurrence. By contrast, fidaxomicin or vancomycin use as first CDI treatment and 5-ASA exposure were associated with negative SHAP values, suggesting a protective association against recurrence.

Additional predictors included a higher Charlson Comorbidity Index, lower hemoglobin, and recent surgery modestly increased recurrence probability, while UC and higher hematocrit or sodium levels were linked to lower predicted risk. Collectively, these SHAP findings provide clinically interpretable insight into

the relative and directional impact of patient comorbidity, IBD subtype, and treatment exposures on CDI recurrence risk.

#### DISCUSSION

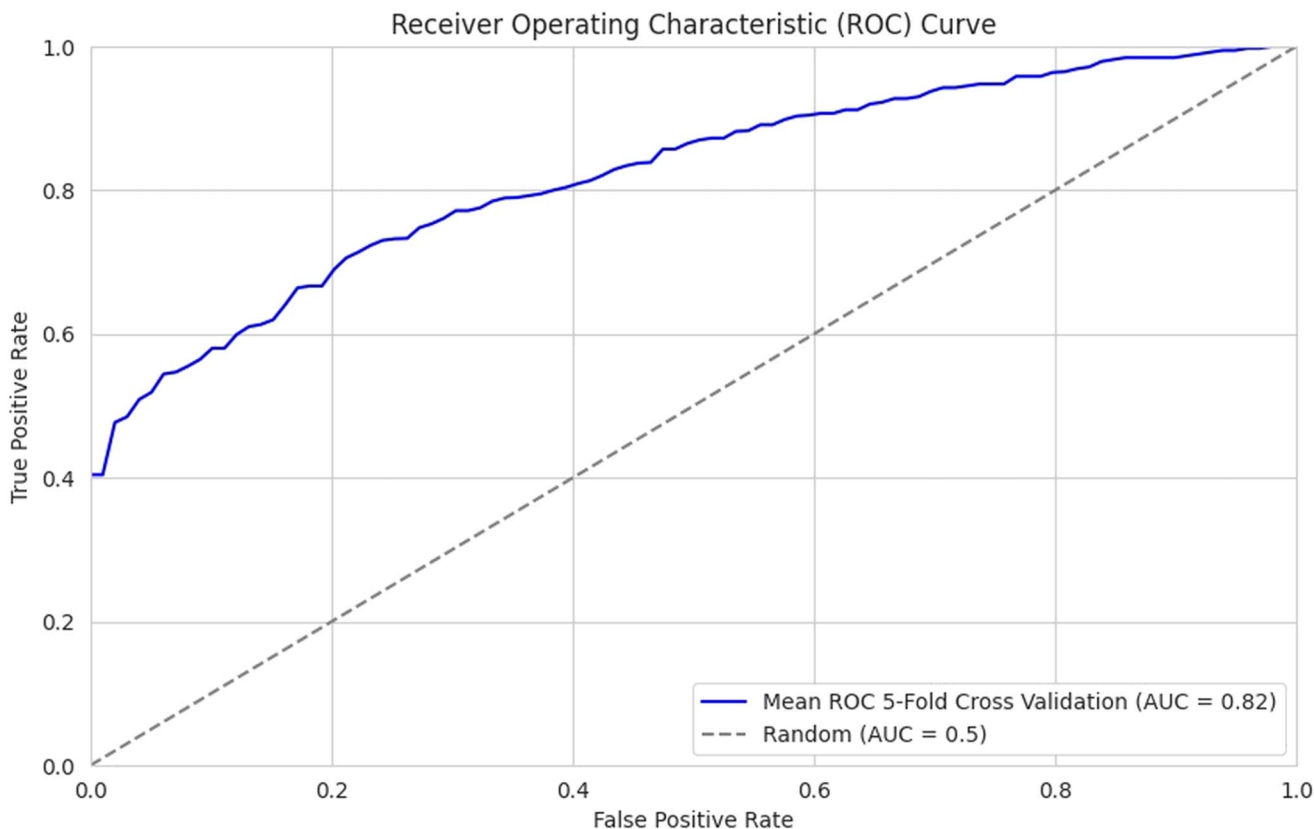
rCDI is a clinically significant complication for patients with IBD, often triggering prolonged illness, escalating therapy, or surgical intervention. Yet identifying which patients are truly at risk remains difficult in day-to-day practice. In this study, we developed a ML model that accurately predicts the risk of CDI recurrence in patients with IBD based on routinely collected clinical data. Our model offers a new tool to proactively identify high-risk patients and enable earlier, more tailored care. From a clinical standpoint, our findings are promising. The model achieved strong performance (AUC of 0.88), with good sensitivity and precision, suggesting it can effectively flag patients who may benefit from closer monitoring or preventive strategies. Importantly, the model draws on information readily available in the

**Table 2.** Mean performance metrics of the Recur CDI-IBD obtained from 5-fold cross-validation on the training data

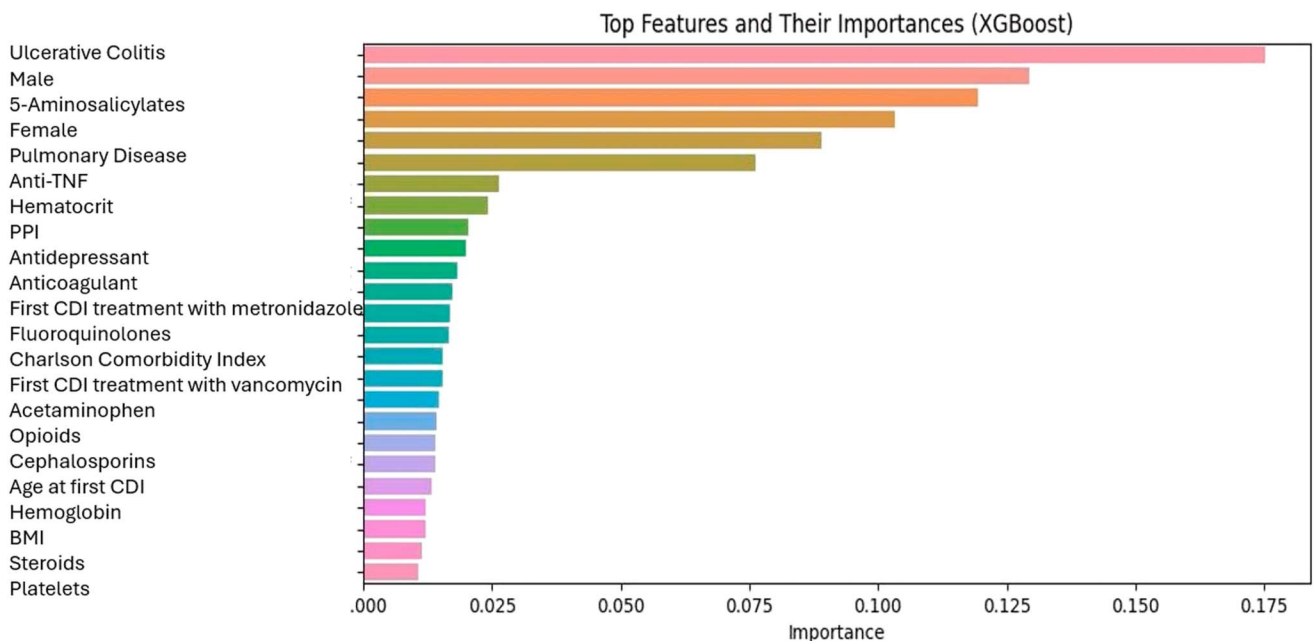
Metric	Mean value	SD ( $\pm$ )
Accuracy	72.75%	2.34%
AUC	0.82	0.03
Sensitivity	0.78	0.09
Specificity	0.75	0.08
Precision	0.72	0.08
Recall	0.71	0.07
F1-score	0.72	0.07

AUC, area under the receiver operating characteristic curve; CDI, *Clostridioides difficile* infection; IBD, inflammatory bowel disease.

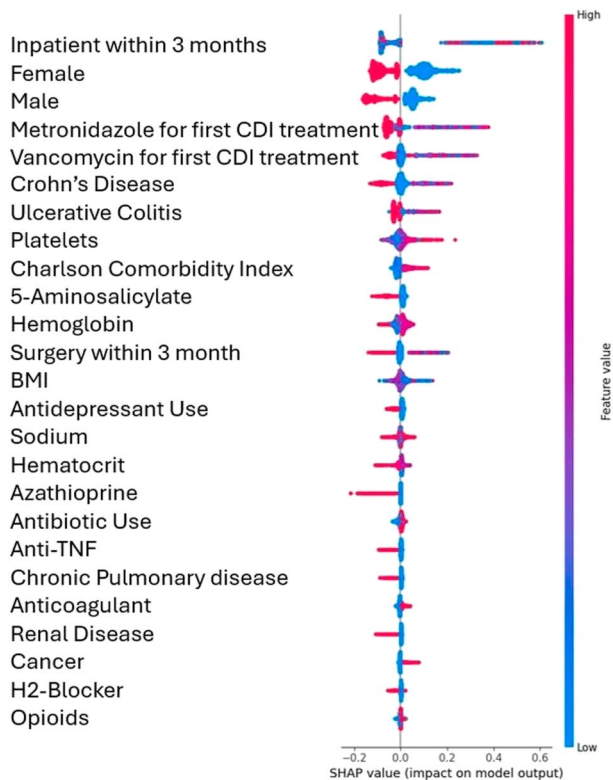
The model demonstrated consistent discrimination between recurrent and nonrecurrent CDI within 60 d.



**Figure 3.** Mean ROC curve obtained during 5-fold cross-validation on the training set (mean AUC = 0.82), confirming the model's internal validity and absence of overfitting. AUC, area under the receiver operating characteristic curve.



**Figure 4.** This plot depicts the top 25 predictors ranked by their relative contribution to the performance of the XGBoost model for predicting recurrent *C. difficile* in patients with inflammatory bowel disease. Each bar represents the normalized importance score derived from the model's gain metric, indicating how much each variable contributed to improving overall model accuracy. BMI, body mass index; CDI, *Clostridioides difficile* infection; PPI, proton pump inhibitor; TNF, tumor necrosis factor.



**Figure 5.** SHAP summary plot illustrates the magnitude and direction of each feature's effect on the model's prediction of recurrent CDI. The x-axis (SHAP value) indicates the feature's contribution to increasing (positive) or decreasing (negative) recurrence probability, while color denotes the feature's raw value (red = higher, blue = lower). Features with higher mean absolute SHAP values exert greater influence on prediction. BMI, body mass index; CDI, *Clostridioides difficile* infection; SHAP, SHapley additive exPlanations; TNF, tumor necrosis factor.

EHR including medications, laboratory values, comorbidities, and recent hospitalizations making it feasible for integration into clinical workflows.

We examined the top 25 predictive features using SHAP and XGBoost importance analysis to enhance the clinical interpretability of our model. Rather than depending on a single biomarker or treatment variable, the model relied on an integrated profile of demographic, disease, treatment, and healthcare utilization data, mirroring how clinicians assess risk in practice. Among the top predictors, UC diagnosis, recent inpatient admission, male sex, corticosteroid or 5-ASA use, and comorbidities including chronic pulmonary and renal disease emerged as consistent high-impact variables. These predictors are clinically meaningful because active disease and recent hospitalizations likely reflect higher inflammatory and antimicrobial exposure burden; corticosteroid use may indicate immune suppression and impaired microbial resilience; and chronic comorbidities may alter host defense and gut microbiome recovery (5). The prominence of UC compared with CD supports existing observations that patients with UC may have greater susceptibility to CDI recurrence, possibly related to colonic microbial dysbiosis and mucosal disruption (4). These findings align with prior literature, a previous study of 137 IBD patients with CDI did not identify any reliable recurrence predictors, (2) whereas another study in 140 Iranian patients with IBD suggested that patients using infliximab in combination with

immunomodulators and steroids had increased risk of incident CDI (18). Conversely, another study reported a protective effect and a lower risk of CDI in patients with UC on anti-TNF biologic therapy (19). From a clinical standpoint, this multifeature approach that our model used, reflects how gastroenterologists already make decisions, by integrating multiple patient factors rather than relying on binary thresholds. It also demonstrates the value of ML in surfacing subtle patterns that might not be captured through traditional regression. (7). A systematic review on the application of Lasso regression in gastroenterology emphasized the importance of including various clinical variables to enhance outcome predictions (20).

To evaluate the predictive performance of our model, we assessed both accuracy and AUC. Accuracy provides a widely used metric that quantifies the proportion of correct classifications made by the model relative to observed clinical outcomes. However, as a single-threshold measure, it does not reflect how well the model distinguishes between patients at high versus low risk across a range of probability thresholds. AUC, in contrast, offers a broader assessment of discriminative ability by quantifying how consistently the model assigns higher predicted risk to patients who actually experience recurrence (21). Taken together, a high AUC paired with strong accuracy suggests that the model is not only correct in its classifications but also reliable in its risk stratification, which is essential for informing personalized treatment decisions in clinical practice.

Our model has several limitations. The data set from which it was derived lacked racial and ethnic diversity and was drawn from a single academic center. Moreover, the relative rarity of CDI recurrence required statistical balancing to ensure model fairness. We did not include imaging data which could have provided complementary insights into IBD severity, complications, and treatment response. Although the Recur CDI-IBD model demonstrated consistent discrimination and balanced predictive performance, its accuracy, precision, and recall should be interpreted with caution. These findings are exploratory and reflect retrospective model development rather than clinical readiness. External validation across independent health systems will be essential to confirm generalizability and determine the model's practical value in guiding patient care.

If validated prospectively, this tool could assist clinicians in early identification of patients with IBD at elevated risk for rCDI, prompting proactive management strategies such as optimizing immunosuppressive regimens, reinforcing infection control practices, or considering early microbiome-targeted therapy. Integration of the RecurCDI-IBD model into existing EHR systems could further enhance its clinical utility by generating automated alerts or risk scores when a patient with IBD is diagnosed with CDI. Such real-time predictions could guide treatment intensity (antibiotic selection and microbiota restoration referral), prompt early consultation with infectious disease and IBD specialists, and enable closer posttreatment monitoring. For successful adoption, future efforts should focus on seamless workflow integration, clear action thresholds, and model transparency to support clinician trust and real-world usability. Over time, this approach could reduce recurrence-related hospitalizations and improve patient quality of life.

In conclusion, our model offers a promising, data-driven means to anticipate rCDI in patients with IBD. By leveraging information already embedded in the EHR, it holds potential to complement clinical judgment, guide early interventions, and improve outcomes

for a population at particularly high risk. Future studies should focus on expanding the data set to include a more diverse patient demographic, ensuring broader applicability. In addition, real-world validation through prospective evaluations will be crucial to further refine and enhance the model's accuracy. Collaborative research efforts, integrating multicenter data, can address potential biases inherent in single-center studies.

#### CONFLICTS OF INTEREST

**Guarantor of the article:** Sahil Khanna, MBBS, MS, FACC.

**Specific author contributions:** R.T.: conceptualization, manuscript writing and revision. A.S.: data analysis. S.A.: data analysis. D.S.P.: supervision, critical revision of manuscript. W.A.F.: supervision, critical revision of manuscript. S.K.: conceptualization, supervision, critical revision of manuscript.

**Financial support:** National Institute of Diabetes and Digestive and Kidney Diseases of the National Institutes of Health under Award Number P30DK084567.

**Potential competing interests:** None to report.

**Ethics statement:** Institutional review board approval was obtained.

**Data availability:** Data, analytic methods, and study materials will be available upon request.

**IRB approval statement:** The study was approved by Mayo Clinic Institutional Review board.

**Availability of data and materials:** All data are available from the corresponding author upon reasonable request.

### Study Highlights

#### WHAT IS KNOWN

- ✓ CDI recurrence is common and clinically burdensome in patients with IBD.
- ✓ Predicting recurrent CDI risk in IBD remains challenging.
- ✓ No validated prediction tool exists for recurrent CDI in IBD populations.

#### WHAT IS NEW HERE

- ✓ A supervised machine-learning model predicts recurrent CDI in IBD using EHR data.
- ✓ The model demonstrates strong discrimination for CDI recurrence risk.
- ✓ SHAP analysis identifies clinically meaningful predictors of recurrent CDI.

#### REFERENCES

1. Khanna S, Shin A, Kelly CP. Management of *Clostridium difficile* infection in inflammatory bowel disease: Expert review from the clinical practice updates committee of the AGA institute. *Clin Gastroenterol Hepatol* 2017;15(2):166–74.
2. Solanky D, Pardi DS, Loftus EV, et al. Colon surgery risk with corticosteroids versus immunomodulators or biologics in inflammatory bowel disease patients with *Clostridium difficile* infection. *Inflamm Bowel Dis* 2019;25(3):610–9.
3. Voth E, Solanky D, Loftus EV Jr, et al. Novel risk factors and outcomes in inflammatory bowel disease patients with *Clostridioides difficile* infection. *Therap Adv Gastroenterol*. 2021;14:1756284821997792. doi: 10.1177/1756284821997792.
4. Nitzan O, Elias M, Chazan B, et al. *Clostridium difficile* and inflammatory bowel disease: Role in pathogenesis and implications in treatment. *World J Gastroenterol* 2013;19(43):7577–85.
5. Khor B, Gardet A, Xavier RJ. Genetics and pathogenesis of inflammatory bowel disease. *Nature* 2011;474(7351):307–17.
6. Ley C, Martin RK, Pareek A, et al. Machine learning and conventional statistics: Making sense of the differences. *Knee Surg Sports Traumatol Arthrosc* 2022; 30(3):753–7.
7. Javaid A, Shahab O, Adorno W, et al. Machine learning predictive outcomes modeling in inflammatory bowel diseases. *Inflamm Bowel Dis* 2022;28(6):819–29.
8. Rajula HSR, Verlatto G, Manchia M, et al. Comparison of conventional statistical methods with machine learning in medicine: Diagnosis, drug development, and treatment. *Medicina (Kaunas)* 2020;56(9):455.
9. Bzdok D, Altman N, Krzywinski M. Statistics versus machine learning. *Nat Methods* 2018;15(4):233–4.
10. Hou JK, Tan M, Stidham RW, et al. Accuracy of diagnostic codes for identifying patients with ulcerative colitis and Crohn's disease in the Veterans Affairs Health Care System. *Dig Dis Sci* 2014;59(10):2406–10.
11. Kelly CR, Fischer M, Allegretti JR, et al. ACG clinical guidelines: Prevention, diagnosis, and treatment of *Clostridioides difficile* infections. *Am J Gastroenterol* 2021;116(6):1124–47.
12. Natekin A, Knoll A. Gradient boosting machines, a tutorial. *Front Neurobot*. 2013;7:21. doi:10.3389/fnbot.2013.00021.
13. Zhang Z. Missing data imputation: Focusing on single imputation. *Ann Transl Med* 2016;4(1):9.
14. Hassanzadeh R, Farhadian M, Rafeemehr H. Hospital mortality prediction in traumatic injuries patients: Comparing different SMOTE-based machine learning algorithms. *BMC Med Res Methodol* 2023;23(1):101.
15. Carrington AM, Manuel DG, Fieguth PW, et al. Deep ROC analysis and AUC as balanced average accuracy, for improved classifier selection, audit and explanation. *IEEE Trans Pattern Anal Mach Intell* 2023; 45(1): 329–41.
16. Brownlee J. Feature importance and feature selection with XGBoost in python. 2020.
17. Rodriguez-Perez R, Bajorath J. Interpretation of machine learning models using shapley values: Application to compound potency and multi-target activity predictions. *J Comput Aided Mol Des* 2020;34(10): 1013–26.
18. Gholam-Mostafaei FS, Yadegar A, Aghdaei HA, et al. Anti-TNF containing regimens may be associated with increased risk of *Clostridioides difficile* infection in patients with underlying inflammatory bowel disease. *Curr Res Transl Med* 2020;68(3):125–30.
19. Ananthakrishnan AN, Oxford EC, Nguyen DD, et al. Genetic risk factors for *Clostridium difficile* infection in ulcerative colitis. *Aliment Pharmacol Ther* 2013;38(5):522–30.
20. Ali H, Shahzad M, Sarfraz S, et al. Application and impact of Lasso regression in gastroenterology: A systematic review. *Indian J Gastroenterol* 2023;42(6):780–90.
21. Ling CX, Huang J, Zhang H. AUC: A better measure than accuracy in comparing learning algorithms. In: *Advances in Artificial Intelligence: 16th Conference of the Canadian Society for Computational Studies of Intelligence, AI 2003, Halifax, Canada, June 11–13, 2003, Proceedings 16 2003*. Springer, 329–41.