

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/138629> holds various files of this Leiden University dissertation.

Author: Zand, A.

Title: Artificial intelligence and e-health for inflammatory bowel diseases: The quest to enhance patient experiences, outcomes and costs

Issue Date: 2020-12-10

CHAPTER 6

Artificial Intelligence for Inflammatory Bowel Diseases (IBD); Developing and Validating Machine Learning Models in Big Data to Predict Negative Outcomes

Submitted

A. Zand^{1,2,3}, Z. Stokes^{1,2}, A. Sharma¹, W.K. van Deen⁴, D.W. Hommes^{1,3}

¹ UCLA Center for Inflammatory Bowel Diseases, Vatche and Tamar Manoukian Division of Digestive Disease, David Geffen School of Medicine, University of California at Los Angeles, CA, USA

² OptumLabs Visiting Fellow, Eden Prairie, MN

³ Leiden University Medical Center, Department of Digestive Diseases, Leiden, the Netherlands

⁴ Cedars-Sinai Center for Outcomes Research and Education, Division of Health Services Research, Cedars-Sinai Medical Center, Los Angeles, CA

Abstract

Background and Aims

The accessibility to Big Data and increased computational resources have paved the way for Artificial Intelligence (AI) to potentially predict adverse health events in complex diseases such as Inflammatory Bowel Diseases (IBD) characterized by considerable heterogeneity and alternating disease states.

Methods

We assessed the feasibility and performance of various statistical and AI models in early prediction of adverse outcomes (hospitalizations, surgeries, long-term steroid and biologics use) for IBD patients using The OptumLabs® Data Warehouse (OLDW), a longitudinal, real-world data asset with de-identified administrative claims and electronic health record (EHR) data, and 108 potentially predictive variables. We built a training model cohort and validated our result in another cohort. We used LASSO and Ridge regressions, Support Vector Machines, Random Forests and Neural Networks and assessed their respective performances and analyzed the strongest predictors to the respective models.

Results

72,178 and 69,165 patients were included in the training and validation set, respectively. In total, 4.1% of patients in the validation set were hospitalized, 2.9% needed IBD-related surgeries, 17% used long term steroids and 13% of patients were initiated with biological therapy. Of the AI models we tested, the Random Forest resulted in the highest accuracy (AUCs 0.71-0.92). The artificial neural network performed well in some but not all of the models (AUCs 0.61-0.90).

Conclusions

This study demonstrates that it is feasible to successfully run complex and novel AI models on large longitudinal data sets of IBD patients (Big Data). These models can be applied for risk stratification and implementation of preemptive measures to avoid adverse outcomes in a clinical setting.

Introduction

The burden of Inflammatory Bowel Disease (IBD) on patients as well as society is large. IBD is a progressive disease with a destructive character and is associated with substantial healthcare costs^{1,2}. Prevention of flares is key to preventing disease progression³⁻⁵. However, the disease course is unpredictable and reliable risk factors for flares are difficult to identify⁵. Finding an approach that identifies patients at risk for disease progression would help to better fine-tune treatment strategies in order to prevent adverse outcomes such as hospitalizations, long term steroid use, the initiation of expensive biologics and surgeries. This could help reduce the substantial costs associated with IBD care and improve long-term outcomes⁶.

The development of healthcare technologies driven by Artificial Intelligence (AI) is expected to see a growth of over \$10 billion in just the next 5 years⁷. With the explosive amount of Electronic Medical Records (EMRs), having doubled in size since 2005, studying patient data is easier now than in any previous era⁸. By taking full advantage of EMR data and, other forms of patient information (e.g. wearables, microbiome/genetic testing, e-health applications, imaging), data driven treatment plans targeted at the disease and individual level could be introduced. The opportunities to construct new strategies and technologies that turn this data into actionable provider recommendations is expected to rapidly grow, as showcased by the immense amount of funding that is going into companies that use AI for healthcare⁹.

Recently, there have been multiple studies that were able to accurately and inexpensively use a subset of AI known as Machine Learning (ML) to predict a variety of outcomes and create distinct classifications for IBD patients (Figure 1)¹⁰⁻¹⁸. Han et al created a gene-based ML classification model to better differentiate between patients with Crohn's disease (CD) and ulcerative colitis (UC)¹⁶. Also using a large sample of genetic data, Wei et al were able to successfully create a genotype-based risk prediction model for IBD¹⁴. Beyond gene-based data, researchers have used AI models with insurance claims data to accurately predict IBD related hospitalization or steroid use within a six-month period¹⁰. This ML approach outperformed more costly biomarker methods of predicting negative outcomes, such as testing for fecal calprotectin. These kind of ML approaches to healthcare have not been limited to IBD¹⁹⁻²³.

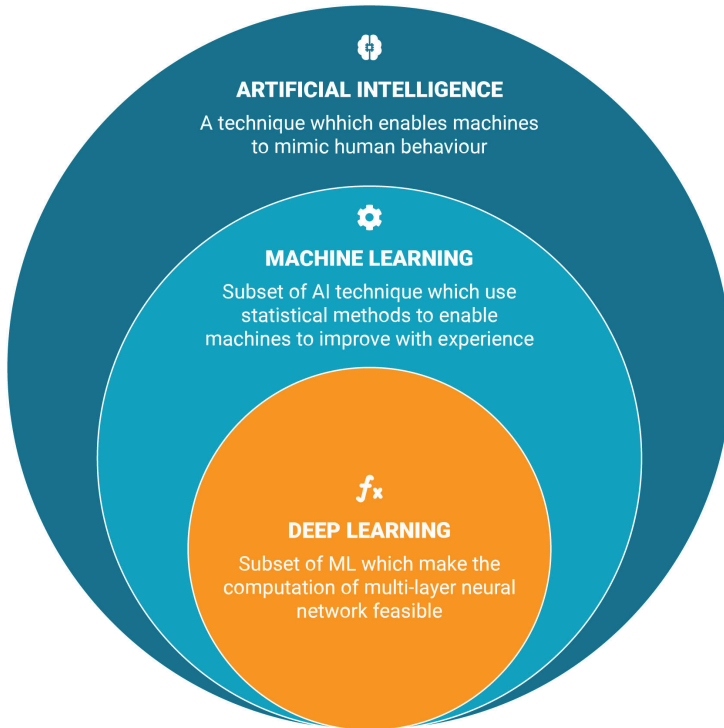


Figure 1. AI is the broad umbrella term of techniques which enables machines to mimic human behavior, when talking about predictive models we usually refer to machine learning which is a subset of AI that uses statistical methods to improve the accuracy of their outcome with experience. Deep Learning is a subset that makes the computation of multi-layer neural networks feasible and thus improving the accuracy even further.

However, studies using the most straightforward data resource, which are administrative databases due to the standardized format and accessibility, to build data driven predictive models for IBD patients were limited in their generalizability. The data came from public health insurance records, while the majority (67.2%) of United States citizens use private insurance, and their samples have limited geographic spread^{13,24}. Additionally, these studies have not attempted to predict other costly negative outcomes such as IBD-related surgeries^{10,13}. To our knowledge, no other study has attempted to apply this ML approach to a larger set of private insurance claims data or use novel deep learning methods such as neural networks. Our goal is to assess the feasibility and performance of various ML models in early prediction of adverse outcomes for IBD patients, including IBD-related surgeries, using a large private insurance claims dataset.

Methods

Study Objectives

The main objective of this study was to assess if variables extracted from insurance claims can predict negative health outcomes in IBD. To achieve this, we assessed the performance of different Machine Learning and Deep Learning models to and compared the performances of the aforementioned models using different performance outcomes.

Data Collection

Deidentified medical, pharmacy and facility claims, were extracted from The OptumLabs[®] Data Warehouse (OLDW), which includes claims from commercially insured individuals and Medicare Advantage beneficiaries (≥ 65 years old) who are representative of the U.S. population with regards to geographical spread, age and race²⁵. Patient-identifying data is removed from the OLDW by OptumLabs before access is granted to investigators. Therefore, this study is not considered human subjects research and is exempt from Institutional Review Board (IRB) regulation.

We created two datasets: a training cohort and a validation cohort. The training cohort contained all patients that were continuously enrolled in their insurance plan between January 1, 2015 and December 31, 2016. The validation cohort includes patients who were continuously enrolled between January 1, 2016 and December 31, 2017. In each cohort, we aimed to predict outcomes in the second year (follow-up) using claims data available in the first year (baseline).

Population

IBD patients were identified using a combination of inpatient and outpatient claims. Patients were included if they had at least two medical claim with diagnosis codes for IBD (International Classification of Diseases, Ninth Revision, Clinical Modification [ICD-9] 555.x or 556.x) **OR** one IBD-related medical claim and one pharmacy claim for IBD-related medication (Supplementary Table 4) in the first year of data.

To ensure enrollees had a specified period of continuous enrollment and the inability to identify an outcome was not due to missing claims data (e.g. enrollee claim was administered by another payor) a continuous enrollment code provided by OLDW was used to make sure the cohorts were continuously enrolled with the respective payor.

Predictive Variables

We constructed 108 variables related to IBD-related care using the claims in the first year of each dataset. These variables were defined based on definitions previously described by

van Deen et al [13]. The variables include the number of IBD-related claims, hospitalizations, emergency department (ED) visits, office visits, procedures, lab and imaging tests, medication use, relapse rate, and comorbidities (for a complete list, see Supplementary Table 1) ¹³.

Model Development

In our models we aimed to predict IBD-related hospitalizations, initiation of biologics, long-term steroid use, and IBD-related surgery in the second year of the data (follow-up) using the 108 utilization-events that occurred in the prior year (baseline). There is consensus in the literature that these are negative outcomes for IBD that should be avoided^{5,6}. *IBD-related hospitalizations* were defined as the presence of any claim for an IBD-related inpatient hospital stay¹³. *Initiation of biologics* was defined as a pharmacy or medical claim for adalimumab, certolizumab pegol, infliximab or natalizumab in the second year, with no claim for that medicine in the first year. *Long-term steroid use* was defined as the use of hydrocortisone, prednisolone, dexamethasone, prednisone and/or methylprednisolone during a consecutive period longer than 90 days based on pharmacy and medical claims. *IBD-related surgery* was defined as any claim with a Current Procedural Terminology (CPT) code specific to an IBD related surgery (See supplementary Table 2 for a full overview).

Logistic and Machine Learning Models

After these datasets were constructed for both cohorts of patients, we trained several logistic regression and machine learning models: a Ridge regression, a LASSO regression, a Support Vector Machine, a Random Forest model, and a Neural Network (See Table 1). Each of these models was trained to predict the probability of a patient incurring a specific negative health outcome in the next year, using the 108 variables from the previous year. We trained five models on the training set of patients and tested them on the validation set.

Ridge regression and LASSO are regression techniques that place a penalty on the model coefficients to ensure that we do not overfit to the training data. Support Vector Machines attempt to separate the patients in the training set who did experience the negative health outcome from those who did not with the largest margin possible. After experimenting with various kernels, we decided on the Gaussian radial basis function. A Random Forest model generates a collection of decision trees, in which each decision tree attempts to find a cut point for each predictor that best separates patients who experienced the negative outcome from those that did not. The cut that achieves the best separation is added to the tree and this process is repeated for each of the two resulting slices of the data, and so on until some minimum number of patients are left in each slice. To capture the nuances in

Table 1. Introduction and Description of Different Models

Model	Explanation	Method	Advantages	Disadvantages
Ridge Logistic	This method creates a model that is not perfectly fit, or overfit, to the data in a given training set. In doing so, it reduces variance and makes the model a better predictor of data points outside of the training set.	Regression	Can reduce overfitting Shrinks effects towards 0 Fast/easy to implement	Simplistic representation may be far from reality Assumptions may be difficult to justify with many predictors
LASSO Logistic	This method attempts to do the same thing as Ridge Regression but uses slightly different mathematical formulas that make it better in certain situations.	Regression	Can reduce overfitting Performs variable selection Fast/easy to implement	Simplistic representation may be far from reality Variable selection is not robust to multicollinearity
Support Vector Machine	Attempts to find the largest separation between two groups. Sometimes the space of observations has to be transformed to find a clear separation.	Machine Learning	Works well with many predictors Makes prediction easy by clearly segmenting population	Lack of a clear separation can lead to poor performance Requires long training times for big data
Random Forest	Random forest is a collection of decision trees trained on different subsets of the data. Each decision tree decides the best places to cut so that observations from the same class fall on the same side of the cut.	Machine Learning	Performs variable selection Good performance for linear and non-linear relationships Fast/easy to implement	Difficult to interpret Prone to overfitting
Neural Network	Neural networks consists of layers of nested linear models (neurons) with a non-linear transformation (activation) after each layer. The output is often the probability that a given observation is a success.	Deep Learning	Captures complex non-linear relationships Fully utilizes big data	Difficult to implement Requires many small decisions that can greatly affect performance

the data, each tree is trained and evaluated on random subsets of the data drawn with replacement. To avoid having too many correlated trees that choose the same best predictors, at each split in the tree only a fraction of the predictors is considered.

Lastly, Neural Networks can identify complex non-linear patterns in the data. These models consist of several imbedded linear functions, known as hidden layers, wrapped in non-linear “activation” functions. These non-linearities in the model work to capture the complicated relationships between the predictors and the probability that a patient will experience the negative outcome. The choice of activation function at each layer plays a big role in determining how well this relationship will be captured by the resulting model. After experimenting with several options, we found that a mix of standard and parametric Rectified Linear Units (ReLUs) performs the best. The last hidden layer is followed by a sigmoid activation function, which outputs a normalized score that we can interpret as the probability that the patient will experience the outcome.

Model Selection Rationale

We trained a battery of machine learning models to discriminate between patients who experienced negative outcomes and those that did not while emphasizing the clinical insights and practical significance that could be understood from the result. To choose the set of base models, on which we would improve with regularization and hyperparameter tuning, we considered the current gap between an algorithm’s complexity/performance and its explainability. We chose several simple linear models with different regularization penalties as they are easy to interpret and align with existing clinical knowledge but often miss complex associations between the variables. We also explored a variety of neural network architectures and tuning procedures to understand the extent to which non-linear relationships in the data could be exploited to improve performance. These models are infamously difficult to understand, as theoretical notions such as statistical significance are difficult to define. With these two extremes covered the SVM and random forest models we considered attempt to strike a balance between performance and interpretability by blending simple structures with complex training procedures. By choosing models that cover this spectrum we can find complicated relationships that lead to solid predictions and warrant prospective validation as well as simpler associations that are easy to validate through expert knowledge.

Performance of the Models

For each model we obtain a prediction for each patient in the validation set. A series of cutoffs were then considered and predictions above the cutoff were labeled as predicted true cases. With these labels the true positive (sensitivity) and true negative (specificity) rates of the model were calculated based on which receiver operating curves (ROC) were constructed. The area under the ROC curve (AUC) for a specific model quantifies the

overall certainty with which the model can predict outcomes at different cut-offs. The single cutoff with the highest geometric average of sensitivity and specificity was selected for each model and specificity and sensitivity values were reported.

Additionally, we calculated the Brier Score which measures the correctness of a model's predictions by summing the differences between the predicted probability of an observation belonging to a class and its actual class label. A low Brier score indicates that the model on average confidently places observations into the correct class. While the AUC quantifies the accuracy of the model, the Brier score quantifies the certainty of the model. For example, if a model assigns a score of 0.51 to every at-risk patient and 0.49 to all other patients, then a cutoff of 0.5 will correctly classify every patient in the validation set and produce a good AUC, but it does not give us a sense of how certain we are about the predictions. The Brier score solves this by measuring the difference between the scores the model predicts (e.g. 0.51) and the true labels (e.g. 1). If all scores are closer to the true label than the Brier score will be close to 0. In this way the Brier score can be used to select the best model from a set with high AUC when the goal is to give not only accurate, but also strong predictions. This is relevant when extrapolating these results to potential meaningful use in a clinical setting.

Feature Importance (except SVM)

The relative importance of the predictive variables in the different models were calculated. For the LASSO and Ridge regression we looked at the magnitude of coefficients and their respective p-values and present the odds ratio. For the Random Forest we measured the importance of each variable by quantifying the change in accuracy of the final predictions after the variable is added to a tree. Larger values indicate the variable is more important. Since the Support Vector Machines did not result in accurate predictions, we did not investigate the relative importance of the predictors. For the neural network we randomly shuffled the observations of a particular variable in the validation set and measured the change in the model's AUC. Variables that create the largest negative change in AUC are defined as the most important.

TRIPOD Statement

Our methodology and research objectives were subject to the TRIPOD (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis) statement which includes a 22-item checklist, which aims to improve the reporting of studies developing, validating, or updating a prediction model, whether for diagnostic or prognostic purposes²⁶. See supplementary table 5 for a full overview.

Tools and Software

Statistical analyses were performed using statistical package program R 3.4.0 and Python.

Results

Population

We included 72,178 patients in our training set and 69,165 patients in our validation set. For both sets the claims from the baseline year (first) were used to generate the 108 predictive features, the follow-up year (second) was used to create our four main outcomes.

Demographics

The mean age of the populations was around 48 years (SD 16.8) for both cohorts and gender was distributed fairly evenly with approximately 52% being female. Both cohorts were predominantly non-Hispanic whites (66% in the training cohort, and 64% in the validation cohort). Looking at medications, biologics use was around 13% for both cohorts in the baseline year, and steroid use was around 27% for both cohorts. We found that 3% of patients in both cohorts had an IBD-related surgery in the baseline year and 6% had an IBD-related hospitalization (Table 2). For a complete overview of the extracted variables during the baseline years of both cohorts, including the average number of hospitalizations, emergency department (ED) visits, insurance coverage, office visits, procedures, lab and imaging tests, and medication use, see Supplementary Table 1.

In the training cohort, 3392 (4.7%) patients had an IBD-related hospitalization, 2454 (3.4%) had IBD-related surgery, 11332 (15.7%) used long term-steroids, and 8661 (12.0%) patients started biological therapy during the one year of follow-up (Table 2).

In the validation cohort, 2863 (4.1%) patients had an IBD-related hospitalization, 2006 (2.9%) had an IBD-related surgery, 11758 (17.0%) used long term steroids, and 9199 (13.3%) of patients started biological therapy during the one year of follow-up (Table 2).

Performance the Validation Model

For the prediction of *IBD-related hospitalizations*, the Random Forest model performed most optimally with an AUC of 0.73 (66% sensitivity, 67% specificity) and a Brier score of 0.21 (See Table 3 and Figure 2). For the prediction of *Initiation of biologics*, the LASSO regression performed best with an AUC of 0.94 (83% sensitivity, 96% specificity) and a Brier Score of 0.05, followed by the Random Forest with an AUC 0.92 (82% Sensitivity, 92% Specificity) and Brier Score of 0.10. Similarly, the Random Forest performed best for

Table 2. Baseline Demographics and Variables of Training and Validation Cohorts in the baseline year

Variable	Training Set Baseline (2015) N= 72,178		Validation Set Baseline (2016) N= 69,165	
Age, mean (SD)	48.5 years (16.8)		47.9 years (16.5)	
Female Gender, n (%)	38254	(53%)	35966	(52%)
Race, n (%)				
White	47710	(66.1%)	44473	(64.3%)
Unknown	12776	(17.7%)	12381	(17.9%)
Black	5052	(7%)	5672	(8.2%)
Hispanic	4692	(6.5%)	4219	(6.1%)
Asian	1949	(2.7%)	2490	(3.6%)
Hospitalizations and ER visits in baseline year, n (%)				
Any ER Visit (#103)	10827	(15%)	11066	(16%)
Any Hospitalization (#97)	4331	(6%)	4150	(6%)
Any IBD-related Hospitalization (#100)	3609	(5%)	3458	(5%)
Any IBD-related ER Visit (#105)	2887	(4%)	2767	(4%)
Any IBD-related surgery (#64)	2165	(3%)	2075	(3%)
Medication use during baseline year, n (%)				
Any IBD Medication use (#1)	28149	(39%)	15908	(23%)
Any Aminosalicylate use (#2&6)	12270	(17%)	11758	(17%)
Any Antibiotic use (#8)	7218	(10%)	6917	(10%)
Any Corticosteroid use (#11,14,17)	18766	(26%)	18675	(27%)
Any Immunomodulator use (#21, 24, 27)	5774	(8%)	5533	(8%)
Any Biologics use (#42)	8661	(12%)	8991	(13%)
Adverse outcomes follow-up year				
	Follow-up year (2016)		Follow-up year (2017)	
IBD-related hospitalizations	3392	(4.70%)	2863	(4.14%)
Initiation of biologics	8661	(12%)	9199	(13.3%)
Long-term steroid Use	11332	(15.7%)	11758	(17%)
IBD-related surgery	2454	(3.4%)	2006	(2.9%)

Refers to the corresponding feature in Supplementary Table 1.

the prediction of *Long-term steroid use* with an AUC of 0.81 (48% Sensitivity, 86% Specificity) and Brier score of 0.15. For the prediction of *IBD-related surgery*, the LASSO Regression and Random Forest had the highest AUC, 0.71 and Brier scores of 0.22 and 0.21, respectively.

Overall, the Random Forest resulted in high AUCs for all outcomes, as did the LASSO regression. The Neural Network performed well for some outcomes, but not others. The Support Vector Machine and Ridge regressions, on the other hand, consistently had lower performance than other models. Of the four outcomes included, the models were able to

predict the initiation of biologics with the highest accuracy, while IBD-related surgery was the most challenging to predict.

Feature Importance

The relative importance of the predictive variables (Supplementary Table 1) in the different models were calculated except the SVM because of its poor performance. To predict *IBD-related hospitalizations*, long-term steroid use and IBD-related surgeries were strong predictors in both the LASSO and Ridge Regressions. Interestingly, the intensity of healthcare utilization as measured by the number of claims or office visits were the strongest predictors in the Random Forest model, which resulted in similar accuracy compared to the regression models. In the Neural Network on the other hand medication use variables were the most important predictors, but with much lower accuracy, indicating that this model was unable to identify the strongest relationship with IBD-related hospitalizations (Table 3).

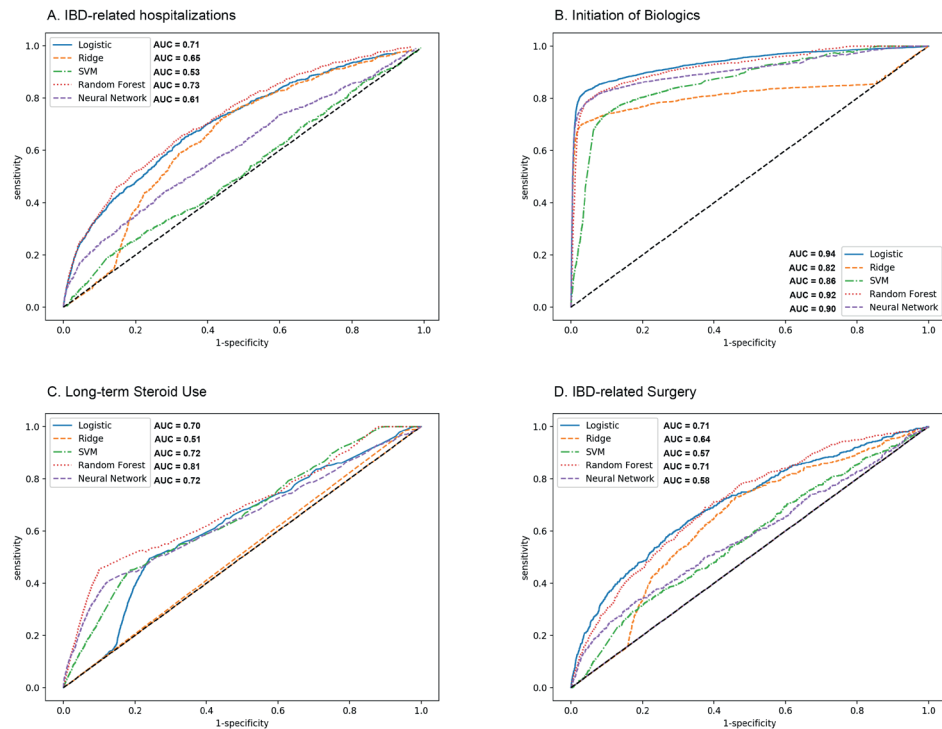


Figure 2. Overview of the performance of the different models for the 4 main outcomes

Table 3. Performance of the different models for the 4 main outcomes

	Sensitivity	Specificity	AUC	Brier Score
IBD-related Hospitalizations				
Ridge Logistic	72%	56%	0.65	0.95
LASSO Logistic	65%	66%	0.71	0.17
Support Vector Machine	54%	48%	0.53	0.04
Random Forest	66%	67%	0.73	0.21
Neural Network	57%	58%	0.61	0.04
Initiation of Biologics				
Ridge Logistic	70%	97%	0.82	0.07
LASSO Logistic	83%	96%	0.94	0.05
Support Vector Machine	75%	89%	0.86	0.10
Random Forest	82%	92%	0.92	0.10
Neural Network	81%	93%	0.90	0.05
Long-term Steroid Use				
Ridge Logistic	99%	4%	0.51	0.83
LASSO Logistic	52%	74%	0.70	0.83
Support Vector Machine	50%	74%	0.72	0.13
Random Forest	48%	86%	0.81	0.15
Neural Network	50%	74%	0.72	0.16
IBD-related surgery				
Ridge Logistic	72%	55%	0.64	0.97
LASSO Logistic	64%	67%	0.71	0.22
Support Vector Machine	54%	55%	0.57	0.03
Random Forest	69%	63%	0.71	0.21
Neural Network	50%	63%	0.58	0.03

Regarding *initiation of biologics*, across all models the use of previous steroids was strongly predictive of a patient being initiated on biologics. The LASSO and Ridge Regressions also found previous CRP lab test and IBD surgeries as strong predictors as well. The random forest, which had the highest accuracy overall, found more heterogeneous predictors including ED visits, number of upper endoscopies and X-ray whereas the neural network mostly found previous use of steroids as the strongest predictor.

Table 4. Feature Importance of the Different Models

The performance of the Support Vector Machine was excluded because of its overall poor performance.

IBD-related Hospitalizations

	Ridge Logistic (AUC = 0.65; Brier score = 0.95)	OR	LASSO Logistic (AUC = 0.71; Brier score = 0.17)	OR	Random Forest (AUC = 0.73; Brier score = 0.21)	Neural Network (AUC = 0.61; Brier score = 0.04)
1	#65 Number of acute IBD surgeries	8.72	#20 Episodes of long-term steroids	1.96	#44 Number of IBD claims	#102 Number of ED visits
2	#64 Any IBD surgeries	2.74	#88 Number of Clostridium difficile stool tests	1.57	#49 Number of office visits	#36 Any certolizumab used this year
3	#88 Number of Clostridium difficile stool tests	2.24	#65 Number of acute IBD surgeries	1.52	#47 Number of UC claims	#35 Episodes of infliximab
4	#20 Episodes of long-term steroids	1.72	#43 Number of episodes of biologics	1.52	#94 Total number of claims	#5 Any oral aminosalicylates used this year
5	#54 Any IBD-related GI visits	1.61	#84 Any MR scans this year	1.51	#96 Number of hospitalizations	#30 Any adalimumab used this year

Initiation of Biologics

	Ridge Logistic (AUC = 0.82; Brier score = 0.07)	OR	LASSO Logistic (AUC = 0.94; Brier score = 0.05)	OR	Random Forest (AUC = 0.92; Brier score = 0.10)	Neural Network (AUC = 0.90; Brier score = 0.05)
1	#42 Any Biologics this year	4.65	#42 Any Biologics this year	8.72	#8 Any antibiotics used this year	#16 Episodes of rectal steroids
2	#13 Episodes of budesonide	2.71	#13 Episodes of budesonide	2.74	#103 Any ED visits this year	#17 Any systemic steroids used
3	#90 Any TB tested this year	2.31	#90 Any TB tested this year	2.24	#10 Episodes of antibiotics	#19 Episodes of systemic steroids
4	#64 Any IBD surgeries	2.29	#23 Episodes of thiopurines	1.72	#80 Any X-rays this year	#20 Episodes of long-term steroids
5	#23 Episodes of thiopurines	2.14	#67 Number of c-reactive protein tests	1.61	#59 Number of upper endoscopies	#21 Any thiopurines used this year

Long-term Steroid Use

	Ridge Logistic (AUC = 0.51; Brier score = 0.83)	OR	LASSO Logistic (AUC = 0.70; Brier score = 0.83)	OR	Random Forest (AUC = 0.81; Brier score = 0.15)	Neural Network (AUC = 0.72; Brier score = 0.16)
1	#20 Episodes of long-term steroids	2.47	#20 Episodes of long-term steroids	2.52	#91 Any influenza vaccine this year	#2 Any rectal aminosalicylates used this year
2	#23 Episodes of thiopurines	2.01	#1 Any IBD medication use	1.61	#103 Any ED visits this year	#7 Episodes of oral aminosalicylates
3	#38 Episodes of certolizumab	1.89	#8 Any antibiotics used this year	1.49	#81 Number of CT scans	#8 Any antibiotics used this year
4	#32 Episodes of adalimumab	1.80	#32 Episodes of adalimumab	1.42	#90 Any TB tested this year	#3 Number of days rectal aminosalicylates used
5	#1 Any IBD medication use	1.58	#78 Any hepatitis B vaccination this year	1.32	#69 Number of sedimentation rate tests	#4 Episodes of rectal aminosalicylates

IBD-related surgery

	Ridge Logistic (AUC = 0.64; Brier score = 0.97)	OR	LASSO Logistic (AUC = 0.71; Brier score = 0.22)	OR	Random Forest (AUC = 0.71; Brier score = 0.21)	Neural Network (AUC = 0.58; Brier score = 0.03)
1	#11 Any budesonide this year	4.85	#108 Any severe disease this year	1.96	#33 Any infliximab used this year	#3 Number of days rectal aminosalicylates used
2	#65 Number of acute IBD surgeries	3.32	#11 Any budesonide this year	1.78	#44 Number of IBD claims	#2 Any rectal aminosalicylates used this year
3	#54 Any IBD-related GI visits	3.18	#65 Number of acute IBD surgeries	1.76	#81 Number of CT scans	#5 Any oral aminosalicylates used this year
4	#84 Any MR scans this year	2.48	#84 Any MR scans this year	1.68	#82 Any CT scans this year	#17 Any systemic steroids used
5	#20 Episodes of long-term steroids	2.48	#20 Episodes of long-term steroids	1.68	#51 Number of IBD office visits	#16 Episodes of rectal steroids

Concerning *long-term steroid use*, the regression models again found previous episodes of IBD medication use to be the strongest predictors. The random forest had the highest accuracy and found medical procedures such as imaging and lab tests and ED visits amongst one of the most predictive features. Similar to initiation of biologics, the neural network found episodes and use of IBD medication, in this particular instance aminosalicylates as the strongest predictor.

Lastly, for our fourth outcome *IBD-related surgery* we found comparable patterns within the regression models showing similar results with episodes of long-term steroids, imaging studies, gastroenterology related visits and severe disease being the greatest predictors. The random forest, which was again one of the best performing models, found infliximab use as the strongest predictor, followed by the total of numbers of IBD-related claims, indicating overall utilization was a strong predictor of IBD-related surgery. Interestingly, the neural net again found use of aminosalicylates as the most predictive feature.

Applying Outcomes in the Daily Clinical Practice

There are several ways that these models can be impactful in daily clinical practice. First, the odds ratios provided by the linear models (ridge logistic and LASSO logistic) can be used to evaluate the risk of patients. For example, we found that risk of hospitalization is strongly linked to previous acute IBD surgeries. Specifically, all else being equal an acute IBD surgery increases the odds of a patient being hospitalized by a factor of more than 8. Second, the complex models that pick up on detailed interactions between the features can be used to make precise risk assessments based on an individual patient's data. As demonstrated by the accuracy of these models, these risk assessments can be used to flag patients that are likely to have a negative outcome with enough notice that providers have time to react and course correct. For example, if we consider a patient with a set of features similar to that of the average patient in the training dataset we can use our models to find that the probability of this patient being hospitalized within the next year is approximately 0.41. This value can give us a sense of the risk assumed by the average IBD patient. Patients whose risk far exceeds this value can be treated as high risk monitored more frequently for predictive markers like CRP or fecal calprotectin.

Lastly, alongside general conclusions about the patient population and risk assessments, these models can be used to evaluate and rank clinical recommendations at the patient level. In this way the models can be used in conjunction with clinical knowledge to motivate actionable, tailored recommendations that are aimed at de-escalating the patient to a lower risk category. Returning to our example of the average patient, we can consider changes to their features that reduce the risk of hospitalization. By examining each feature individually,

the model finds that similar patients to this one benefit from a *Clostridium difficile* stool test. Specifically, our patient is forecasted to see a reduction in their probability of being hospitalized from 0.41 to approximately 0.29 as a result of this intervention. Between these three applications of our results to clinical practice it is clear that the models we have found provide the foundation for a novel, targeted approach to data-driven IBD care.

Discussion

This study demonstrated that it was feasible to successfully run complex machine learning models on large (Big Data) and representative longitudinal claims data sets of IBD patients. We analyzed traditional models including LASSO and Ridge regressions, machine learning methods such as Support Vector Machines and Random Forests but also included more novel methods like Neural Networks, and successfully compared their relative performance. Overall, the Random Forest performed best across all outcomes, which might indicate that the relationships between the claim's features are best captured by a Random Forest model and that this model framework might work best for claims predictions in general.

Regarding feature importance, it is worth noting that the models returned different features for the different outcomes. The regression models overall had comparable findings, with the most predictive features of negative outcomes being largely related to medication use. The random forest had the highest accuracy overall but had more heterogenous findings, being less limited to medication use as the most predictive feature but also including procedures such as imaging and lab tests as strong predictors. Lastly, the neural net had the most consistent findings across all outcomes, which were mostly medication use related. The difference in findings across the models would argue for the need to explore various models depending on the available data and the choice of outcomes. Based on the research objectives and available data, the models can expose different outcomes and relationships, and this can have an impact on the interpretation and clinical implementation. Furthermore, more novel methods such as neural networks should be further investigated and explored in order to increase accuracy and to examine if they can potentially expose correlations and non-linear relationships that might not be found in more conventional methods.

Several others have used claims data to predict IBD-related utilization events in specific IBD sub-populations. For instance, Waljee et al. applied their model to a set of Veteran's Health Administration data, which limited their sample to a 93% male and old (mean age

59 years) population¹⁰; furthermore, public insurance is only used by a minority of United States population²⁴. Other prior works that have used ML approaches on private insurance data have been limited by the geographic spread of their sample¹³. To our knowledge, this is the first study utilize this ML based prediction approach on a nationally representative IBD population. Additionally, different outcomes were used in some of these studies. Waljee et al. used a composite measure capturing both hospitalization and corticosteroid use, where we have split up these outcomes and checked for long-term steroid use. Their composite measure had an AUC of 0.85 and Brier score of 0.20. We found similar results in our Random Forest model with a AUC of 0.73 and Brier score of 0.21 for hospitalizations and 0.81 AUC and 0.15 Brier Score for long-term steroid use. Furthermore, to our knowledge, our study is the first to predict IBD-related surgery using claims data. Additionally, to our knowledge, the use of novel deep learning methods such as Neural Networks has not been described previously in the IBD literature. These new methods should be further explored and reported on as they have the potential to unlock new opportunities for personalized management in IBD and also because of the fact that these models are now feasible to run because of the increased availability of Big Data and increased computational resources.

There are some limitations worth noting to this study. While a data driven approach to healthcare has great potential to improve patient outcomes, there are some limitations to ML that are worth noting. For one, ML algorithms can only describe correlations between variables or features of interest, not necessarily causation²⁷. Furthermore, assumptions are generally made about data sets when applying a given ML algorithm to it, which can narrow the scope of the model in real world situations²⁷. In our case, we pre-defined 108 variables to include in our model. Additionally, some outcomes may have a more complicated (i.e. non-linear) relationship with the predictors, and the models we chose may not capture those relationships. Also, we did not include data from the EMR in our prediction model, inclusion of clinical variables could improve the predictive accuracy. However, administrative databases are more readily accessible due to the standardized format and are therefore remain a more straightforward source of data for these initiatives.

Looking ahead, the practical reality of AI is an enigma to many practitioners (See Figure 1 and Table 1). With boundless publications discussing the new wealth of electronic databases and promises of “Big Data”, most never go into details about what exactly these new technologies are doing to, for example, “outperform cardiologists reading EKGs”²⁹. Unlike the days of small data sets collected through calculated experiment and observation,

this data cannot be studied with the standard methods of statistical analysis⁹. The computations that are generally feasible in experimental settings require vast computational resources when the data is on the order of millions of observations. Therefore, smarter algorithms were created to perform statistical analysis on large data sets. Many would refer to this jump as the development of Machine Learning (ML), but formally it is closer to the sub-field of Computational Statistics. The real jump to ML utilizes the vast amounts of data in a sophisticated way that emphasizes accurate predictions of outcomes over significance and interpretability⁹. With this mindset change, outcomes can be evaluated by experts and the entire process can be incorporated into decision support in daily clinical practice. Now, without much effort from the user, algorithms can make predictions given new data and automatically make a recommendation or perform some action, appearing to have Artificial Intelligence (AI)⁹. With the increase of computational power and abundance of longitudinal patient data, applying machine learning and its subset of Deep Learning in Big Data sets has become feasible. In this study we provide the first steps in this direction. Kim et al. (2019) has already showcased transferability of these models to different institutions, alleviating a major concern¹⁹. The next step would be to integrate these models in a prospective setting to study their performance on reliability, patient outcomes and costs.

References

1. Pariente B, Cosnes J, Danese S, et al. Development of the Crohn's disease digestive damage score, the Lémann score. *Inflammatory Bowel Diseases*. 2011;17(6):1415-1422. doi:10.1002/ibd.21506
2. Kappelman MD, Rifas-Shiman SL, Porter CQ, et al. Direct Health Care Costs of Crohn's Disease and Ulcerative Colitis in US Children and Adults. *Gastroenterology*. 2008;135(6):1907-1913. doi:10.1053/j.gastro.2008.09.012
3. D'Haens G, Baert F, van Assche G, et al. Early combined immunosuppression or conventional management in patients with newly diagnosed Crohn's disease: an open randomised trial. *The Lancet*. 2008;371(9613):660-667. doi:10.1016/S0140-6736(08)60304-9
4. Kang B, Choi SY, Kim HS, Kim K, Lee YM, Choe YH. Mucosal healing in paediatric patients with moderate-to-severe luminal Crohn's disease under combined immunosuppression: Escalation versus early treatment. *Journal of Crohn's and Colitis*. Published online 2016. doi:10.1093/ecco-jcc/jjw086
5. Olivera P, Danese S, Jay N, Natoli G, Peyrin-Biroulet L. Big data in IBD: a look into the future. *Nature Reviews Gastroenterology and Hepatology*. 2019;16(5):312-321. doi:10.1038/s41575-019-0102-5
6. van der Valk ME, Mangen MJJ, Severs M, et al. Evolution of costs of inflammatory bowel disease over two years of follow-up. *PLoS ONE*. 2016;11(4). doi:10.1371/journal.pone.0142481
7. Statista. Global AI software market size 2018-2025 | Statista. Tractica. Published 2019. Accessed July 20, 2020. <https://www.statista.com/statistics/607716/worldwide-artificial-intelligence-market-revenues/>
8. Office-based Physician Electronic Health Record Adoption. Accessed June 24, 2020. <https://dashboard.healthit.gov/quickstats/pages/physician-ehr-adoption-trends.php>
9. Derrington D. *Artificial Intelligence for Health and Health Care.*; 2017. Accessed June 24, 2020. <https://pdfs.semanticscholar.org/4f32/7be94508a5c1f2a6f09917d7dcf57698af24.pdf>
10. Waljee AK, Lipson R, Wiitala WL, et al. Predicting Hospitalization and Outpatient Corticosteroid Use in Inflammatory Bowel Disease Patients Using Machine Learning. *Inflammatory Bowel Diseases*. 2018;24(1):45-53. doi:10.1093/ibd/izz007
11. AK W, B L, K S, et al. Predicting Corticosteroid-Free Biologic Remission with Vedolizumab in Crohn's Disease. *Inflammatory bowel diseases*. 2018;24(6). doi:10.1093/IBD/IZY031
12. Waljee AK, Wallace BI, Cohen-Mekelburg S, et al. Development and Validation of Machine Learning Models in Prediction of Remission in Patients With Moderate to Severe Crohn Disease. *JAMA network open*. 2019;2(5):e193721. doi:10.1001/jamanetworkopen.2019.3721
13. Vaughn DA, van Deen WK, Kerr WT, et al. Using insurance claims to predict and improve hospitalizations and biologics use in members with inflammatory bowel diseases. *Journal of Biomedical Informatics*. 2018;81:93-101. doi:10.1016/j.jbi.2018.03.015
14. Wei Z, Wang W, Bradfield J, et al. Large sample size, wide variant spectrum, and advanced machine-learning technique boost risk prediction for inflammatory bowel disease. *American Journal of Human Genetics*. 2013;92(6):1008-1012. doi:10.1016/j.ajhg.2013.05.002
15. Menti E, Lanera C, Lorenzoni G, et al. Bayesian Machine Learning Techniques for revealing complex interactions among genetic and clinical factors in association with extra-intestinal Manifestations in IBD

- patients. *AMIA . Annual Symposium proceedings AMIA Symposium*. 2016;2016:884-893. Accessed July 27, 2020. /pmc/articles/PMC5333221/?report=abstract
16. Han L, Maciejewski M, Brockel C, et al. A probabilistic pathway score (PROPS) for classification with applications to inflammatory bowel disease. *Bioinformatics*. 2018;34(6):985-993. doi:10.1093/bioinformatics/btx651
 17. Cai T, Lin TC, Bond A, et al. The association between arthralgia and vedolizumab using natural language processing. *Inflammatory Bowel Diseases*. 2018;24(10):2242-2246. doi:10.1093/ibd/izy127
 18. Hou JK, Chang M, Nguyen T, et al. Automated identification of surveillance colonoscopy in inflammatory bowel disease using natural language processing. *Digestive Diseases and Sciences*. 2013;58(4):936-941. doi:10.1007/s10620-012-2433-8
 19. Kim E, Carballo PJ, Castro MR, Pieczkiewicz DS, Simon GJ. Towards more Accessible Precision Medicine: Building a more Transferable Machine Learning Model to Support Prognostic Decisions for Micro- and Macrovascular Complications of Type 2 Diabetes Mellitus. *Journal of Medical Systems*. 2019;43(7). doi:10.1007/s10916-019-1321-6
 20. Nori VS, Hane CA, Martin DC, Kravetz AD, Sanghavi DM. Identifying incident dementia by applying machine learning to a very large administrative claims dataset. *PLoS ONE*. 2019;14(7). doi:10.1371/journal.pone.0203246
 21. Chen S, Bergman D, Miller K, Kavanagh A, Frownfelter J, Showalter J. Using applied machine learning to predict healthcare utilization based on socioeconomic determinants of care. *The American journal of managed care*. 2020;26(1):26-31. doi:10.37765/ajmc.2020.42142
 22. Xiao J, Ding R, Xu X, et al. Comparison and development of machine learning tools in the prediction of chronic kidney disease progression. *Journal of Translational Medicine*. 2019;17(1). doi:10.1186/s12967-019-1860-0
 23. Chiu YC, Chen HIH, Zhang T, et al. Predicting drug response of tumors from integrated genomic profiles by deep neural networks. *BMC Medical Genomics*. 2019;12(Suppl 1). doi:10.1186/s12920-018-0460-9
 24. Kinney ED. Health Insurance Coverage in the United States. In: *Protecting American Health Care Consumers*.; 2020:23-40. doi:10.2307/j.ctv11smv14.6
 25. OptumLabs and OptumLabs Data Warehouse (OLDW) Descriptions and Citation. Cambridge, MA: n.p., May 2019. PDF Reproduced with permission from OptumLabs.
 26. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The TRIPOD statement. *Annals of Internal Medicine*. 2015;162(1):55-63. doi:10.7326/M14-0697
 27. The Limitations of Machine Learning | by Matthew Stewart, PhD Researcher | Towards Data Science. Accessed July 27, 2020. <https://towardsdatascience.com/the-limitations-of-machine-learning-a00e0c3040c6>

Supplementary Table 1. 110 predictive features

Summary of the overall prevalence of the 110 potentially predictive factors included in our models. These features were compiled by experts in the IBD field (WD and DH) and pulled from the 110 features earlier published by WD and DH (Vaughn et al. 2018). Shading of pink represents binary variables, yellow represents variables related to days, and cyan represents variables related to courses of medication. Values reflect the MEAN, which included claims submitted to United HealthCare between 2015 and 2017. Two of the features could not be constructed and were excluded from the analysis (#55 and #93)

#	explanation	Training 2015 mean	Validation 2016 mean	Comparison Vaughn et al.
1	Any IBD related medications use (all the medications in variables #2 - #41)	0.39	0.23	0.88
2	Any rectal aminosalicylates used in this year	0.03	0.03	0.14
3	Number of days rectal aminosalicylates used	99.51	100.39	15
4	Number of times an episode of rectal aminosalicylates started	0.03	0.04	0.17
5	Any oral aminosalicylates used in this year	0.14	0.14	0.53
6	Number of days oral aminosalicylates used	27.78	28.73	124
7	Number of times an episode of oral aminosalicylates started	0.15	0.15	0.47
8	Any antibiotics used in this year	0.1	0.1	0.24
9	Number of days antibiotics used	2.03	1.9	6.6
10	Number of times an episode of antibiotics started	0.1	0.1	0.32
11	Any budesonide (local release steroid) used in this year	0.04	0.04	0.06
12	Number of days budesonide (local release steroid) used	3.43	3.62	7.7
13	Number of times an episode of budesonide (local release steroid) started	0.04	0.04	0.07
14	Any rectal steroids used in this year	0.09	0.09	0.08
15	Number of days rectal steroids used	2.11	2.1	3.9
16	Number of times an episode of rectal steroids started	0.09	0.09	0.10
17	Any systemic steroids used in this year	0.13	0.14	0.28
18	Number of days systemic steroids used	4.87	4.98	19
19	Number of times an episode of systemic steroids started	0.13	0.14	0.39
20	Number of times an episode of long term (>3 consecutive months) steroids started	0.15	0.17	0.06
21	Any thiopurines used in this year	0.06	0.06	0.19
22	Number of days thiopurines used	12.1	12.5	48
23	Number of times an episode of thiopurines started	0.06	0.06	0.13
24	Any methotrexate used in this year	0.01	0.01	0.03
25	Number of days methotrexate used	1.63	1.95	6.0
26	Number of times an episode of methotrexate started	0.01	0.01	0.03
27	Any cyclosporine or tacrolimus used in this year	0.01	0.01	0.01
28	Number of days on cyclosporine or tacrolimus	1.02	0.99	1.4

Supplementary Table 1. Continued

#	explanation	Training 2015 mean	Validation 2016 mean	Comparison Vaughn et al.
29	Number of times an episode of cyclosporine or tacrolimus was started	0.01	0.01	0.00
30	Any adalimumab used in this year	0.03	0.04	0.06
31	Number of days adalimumab used	6.22	8.1	15
32	Number of times an episode of adalimumab started	0.03	0.04	0.04
33	Any infliximab used in this year	0.08	0.09	0.11
34	Number of days infliximab used	20.47	22.13	28
35	Number of times an episode of infliximab started	0.08	0.09	0.08
36	Any certolizumab used in this year	0.01	0.01	0.01
37	Number of days certolizumab used	1.82	1.43	2.9
38	Number of times an episode of certolizumab started	0.01	0.01	0.01
39	Any natalizumab used in this year	0	0	0.00
40	Number of days natalizumab used	0.1	0.11	0.40
41	Number of times an episode of natalizumab started	0	0	0.00
42	Any biologics (variables #30-#41) used in this year	0.12	0.13	0.18
43	Number of times an episode of biologics (variables #30-#41) started	0.12	0.13	0.13
44	Number of IBD claims	20.45	23.21	5.9
45	Number of Crohn's disease claims	12.58	14	3.3
46	Any Crohn's disease claims this year	0.41	0.42	0.51
47	Number of ulcerative colitis claims	8.04	9.41	2.7
48	Any ulcerative colitis claims this year	0.47	0.51	0.63
49	Number of office visits	8.47	8.39	8.1
50	Any office visits this year	0.96	0.96	0.98
51	Number of IBD related office visits	1.73	1.87	2.3
52	Any IBD related office visits this year	0.62	0.65	0.80
53	Number of IBD related office visits with a gastroenterologist	0	0	1.2
54	Any IBD related office visits with a gastroenterologist this year	0	0	0.53
55	Number of IBD related office visits with a UCLA gastroenterologist	N/A	N/A	0.02
56	Any IBD related office visits with a non-UCLA gastroenterologist this year	0.02	0.01	0.51
57	Number of colonoscopies	0.39	0.41	0.49
58	Any colonoscopies this year	0.32	0.34	0.44
59	Number of upper endoscopies	0.11	0.11	0.14
60	Any upper endoscopies this year	0.11	0.11	0.13
61	Number of endoscopies of the small intestine	0	0	0.03

Supplementary Table 1. Continued

#	explanation	Training 2015 mean	Validation 2016 mean	Comparison Vaughn et al.
62	Any endoscopies of the small intestine this year	0	0	0.02
63	Number of IBD related surgeries	0.05	0.06	0.06
64	Any IBD related surgeries this year	0.03	0.03	0.04
65	Number of acute IBD related surgeries (this is a subset of IBD related surgeries)	0.05	0.05	0.06
66	Any acute IBD related surgeries (this is a subset of IBD related surgeries) this year	0.03	0.03	0.04
67	Number of C-reactive protein tests	0.27	0.29	0.68
68	Any C-reactive protein tests this year	0.27	0.29	0.32
69	Number of sedimentation rate tests	0.27	0.28	0.89
70	Any sedimentation rate tests this year	0.25	0.26	0.39
71	Number of stool calprotectin tests	0.04	0.05	0.03
72	Any stool calprotectin tests this year	0.04	0.05	0.02
73	Number of complete blood counts	1.02	1.04	2.7
74	Any complete blood counts this year	0.76	0.77	0.82
75	Number of liver enzyme tests	1.01	1.04	2.3
76	Any liver enzyme tests this year	0.73	0.75	0.79
77	Number of Hepatitis B tests	0.23	0.26	0.12
78	Any hepatitis B vaccination this year	0.1	0.11	0.10
79	Number of X-rays	0.15	0.14	0.24
80	Any X-rays this year	0.11	0.1	0.13
81	Number of CT scans	0.23	0.23	0.29
82	Any CT scans this year	0.20	0.19	0.19
83	Number of MR scans	0.08	0.08	0.06
84	Any MR scans this year	0.05	0.06	0.05
85	Number of ultrasounds	0.08	0.08	0.10
86	Any ultrasounds this year	0.07	0.07	0.08
87	Any bone loss assessment this year	0.06	0.06	0.07
88	Number of Clostridium difficile stool tests	0.1	0.09	0.16
89	Any Clostridium difficile stool tests this year	0.09	0.09	0.12
90	Any TB tested this year	0.09	0.11	0.08
91	Any influenza vaccine this year	0.16	0.13	0.17
92	Any pneumococcal vaccine this year	0.06	0.06	0.02
93	Charlson comorbidity score (higher score implies comorbidities)	N/A	N/A	0.51

Supplementary Table 1. Continued

#	explanation	Training 2015 mean	Validation 2016 mean	Comparison Vaughn et al.
94	Total number of claims	27.45	27.82	73
95	Total number of days prescriptions were covered by plan	113.59	119.61	364
96	Number of hospitalizations	0.13	0.11	0.28
97	Any hospitalizations this year	0.06	0.06	0.17
98	Total number of days hospitalized	0.85	0.74	1.8
99	Number of IBD related hospitalizations	0.16	0.16	0.10
100	Any IBD related hospitalizations this year	0.05	0.05	0.08
101	Total number of days hospitalized related to IBD	0.52	0.53	0.76
102	Number of ED visits	0.72	0.73	0.58
103	Any ED visits this year	0.15	0.16	0.26
104	Number of IBD related ED visits	0.06	0.07	0.24
105	Any IBD related ED visits this year	0.04	0.04	0.13
106	Age	50.13	48.68	42
107	Any moderate disease this year (based on a combination of number of relapses and long term steroid use)	0.01	0	0.21
108	Any severe disease this year (based on a combination of number of relapses and long term steroid use)	0.13	0.19	0.15
109	Relapse rate (based on how use of systemic steroids, use of biologics, and acute IBD related surgeries)	0.06	0.06	0.58
110	The number of years someone has been a continuous member of United HealthCare or Anthem	1.96	2.56	1.6

Supplementary Table 2. Development of Main Outcomes

Hospitalization	For each patient take all claims with place of service code = 21 (inpatient hospital). Next check the 9 diagnosis codes for each hospital claim for any of the following IBD-related ICD 9/10 codes: 5551, 5552, 5559, 5561, 5562, 5563, 5564, 5565, 5566, 5568, 5569, K500, K501, K508, K509, K510, K512, K513, K514, K515, K518, K519. If any of these codes are present in any of the hospitalization claims, then the patient is considered to have had an IBD-related hospitalization that year.
Biologics	For each patient search for facility and pharmacy claims with any of the following drug names or CPT codes: ADALIMUMAB, CERTOLIZUMAB PEGOL, INFlixIMAB, NATALIZUMAB, J0135, J1745, J2323, Q4079, J0718, C9249. If any claims are found, then the patient is considered to have initiated Biologics that year.
Surgery	For each patient search the medical and facility claims for any claims with the following CPT codes: 44005-44346, 44602-44701, 45000-45190, 45395-45999, 46020-46060, 46270-46288, 49000-49084. If any claims are found, then the patient is considered to have had an IBD-related surgery that year.
Long-term Steroids	For each patient search for claims where any of the following steroids were given: HYDROCORTISONE, PREDNISOLONE, DEXAMETHASONE, PREDNISONE, METHYLPREDNISOLONE. Using the variable COUNT_DAYS_SUPPLY calculate the length of time of each episode of steroids. If any episode lasts longer than 3 months (90 days), then the patient is considered to have had an episode of long-term steroids for that year.

Supplementary Table 3. Technical Appendix Models

Model	Technical Detail
Ridge regression and LASSO	The first two models fit include Ridge regression and LASSO. These are regression techniques that place a penalty on the model coefficients to ensure that we do not overfit to the training data. In this way these methods jointly perform variable selection and model training. The primary difference between these two models is in the choice of penalty. Ridge regression penalizes the sum of squares of the least squares estimates and as the user-selected size of the penalty increases all estimates become increasingly smaller but never reach 0. This can be problematic for researchers who are interested in the substantive interpretation of all coefficients in the model. LASSO corrects this problem by instead penalizing the sum of absolute values of the estimates. This change leads some of the estimates to become 0 as the size of the penalty increases. The resulting model then consists only of the estimates that are significantly large.
Support Vector Machine	We also trained several Support Vector Machines with varying kernels. These models attempt to separate the patients in the training set who did experience the negative health outcome from those who did not with the largest margin possible. Since many high-dimensional data are not separable with linear support vectors, transformations through the use of kernels are employed to achieve non-linear regions. We try several such kernels, but the one which obtains the highest testing accuracy, which also happens to be one of the most often used kernels, is the Gaussian radial basis function.
Random Forest	To isolate important variables, we also fit Random Forest models. These are ensemble classifiers made up of collections of decision trees. Each decision tree makes linear cuts through the variable space to achieve the best division between the two classes. To capture the nuances in the data each tree is trained and evaluated on random subsets of the data drawn with replacement. To avoid having too many correlated trees that choose the same best predictors, at each split in the tree only a fraction of the predictors is considered.
Neural Networks	To understand the complex non-linear patterns in the data, we train several Neural Networks. These models consist of several imbedded linear functions, known as hidden layers, wrapped in non-linear “activation” functions. The choice of activation function at each layer determines the functional form of the model. After experimenting with several options we use a mix of standard and parametric Rectified Linear Units (ReLUs). The output layer is followed by a sigmoid activation function, so that we may interpret the output as the probability that the patient will experience the outcome. To train the model we use stochastic gradient descent to minimize a binary cross entropy loss.

Supplementary Table 4. Medications

Drug group	Drug type	Included drugs	CPT
Aminosalicylates	ASA - oral ASA - rectal	mesalamine, sulfasalazine, balsalazide, olsalazine	
Antibiotics		metronidazole, ciprofloxacin	
Corticosteroids	budesonide systemic rectal	budesonide prednisone, methylprednisolone, hydrocortisone, prednisolone, dexamethasone	
Immunomodulators	thiopurines methotrexate cyclosporine tacrolimus	azathioprine, mercaptopurine, methotrexate cyclosporine tacrolimus	
Biologics	adalimumab certolizumab infliximab natalizumab	adalimumab certolizumab pegol infliximab natalizumab	J0135 J0718, C9294 J1745 J2323, Q4079

Supplementary Table 5. TRIPOD Checklist: Prediction Model Development

Section/Topic	Item	Checklist Item	Page
Title and abstract			
Title	1	Identify the study as developing and/or validating a multivariable prediction model, the target population, and the outcome to be predicted.	1
Abstract	2	Provide a summary of objectives, study design, setting, participants, sample size, predictors, outcome, statistical analysis, results, and conclusions.	2
Introduction			
Background and objectives	3a	Explain the medical context (including whether diagnostic or prognostic) and rationale for developing or validating the multivariable prediction model, including references to existing models.	4,5
	3b	Specify the objectives, including whether the study describes the development or validation of the model or both.	4,5
Methods			
Source of data	4a	Describe the study design or source of data (e.g., randomized trial, cohort, or registry data), separately for the development and validation data sets, if applicable.	6
	4b	Specify the key study dates, including start of accrual; end of accrual; and, if applicable, end of follow-up.	6
Participants	5a	Specify key elements of the study setting (e.g., primary care, secondary care, general population) including number and location of centres.	6
	5b	Describe eligibility criteria for participants.	6
	5c	Give details of treatments received, if relevant.	6
Outcome	6a	Clearly define the outcome that is predicted by the prediction model, including how and when assessed.	7
	6b	Report any actions to blind assessment of the outcome to be predicted.	7
Predictors	7a	Clearly define all predictors used in developing or validating the multivariable prediction model, including how and when they were measured.	7
	7b	Report any actions to blind assessment of predictors for the outcome and other predictors.	7
Sample size	8	Explain how the study size was arrived at.	6,7
Missing data	9	Describe how missing data were handled (e.g., complete-case analysis, single imputation, multiple imputation) with details of any imputation method.	6,7
Statistical analysis methods	10a	Describe how predictors were handled in the analyses.	6
	10b	Specify type of model, all model-building procedures (including any predictor selection), and method for internal validation.	7,8
	10d	Specify all measures used to assess model performance and, if relevant, to compare multiple models.	9
Risk groups	11	Provide details on how risk groups were created, if done.	7

Supplementary Table 5. Continued

Section/Topic	Item	Checklist Item	Page
Results			
Participants	13a	Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful.	11
	13b	Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome.	11
Model development	14a	Specify the number of participants and outcome events in each analysis.	11
	14b	If done, report the unadjusted association between each candidate predictor and outcome.	11
Model specification	15a	Present the full prediction model to allow predictions for individuals (i.e., all regression coefficients, and model intercept or baseline survival at a given time point).	11,12
	15b	Explain how to use the prediction model.	13
Model performance	16	Report performance measures (with CIs) for the prediction model.	11,12
Discussion			
Limitations	18	Discuss any limitations of the study (such as nonrepresentative sample, few events per predictor, missing data).	16
Interpretation	19b	Give an overall interpretation of the results, considering objectives, limitations, and results from similar studies, and other relevant evidence.	15
Implications	20	Discuss the potential clinical use of the model and implications for future research.	17
Other information			
Supplementary information	21	Provide information about the availability of supplementary resources, such as study protocol, Web calculator, and data sets.	18
Funding	22	Give the source of funding and the role of the funders for the present study.	1

