# Machine Learning Based Prediction of Surgical Outcomes in Chronic Rhinosinusitis from Clinical Data

Sayeed Shafayet Chowdhury[a,c], Karen D'Souza[b], V. Siva Kakumani[a], Snehasis Mukhopadhyay[a], Shiaofen Fang[c], Rodney J. Schlosser[d], Daniel M. Beswick[e], Jeremiah A. Alt[f], Jess C. Mace[g], Zachary M. Soler[d], Timothy L. Smith[g], Vijay R. Ramakrishnan[h]

[a]*Purdue University, Indianapolis, IN, USA*
[b]*Idaho National Laboratory, Idaho Falls, ID, USA*
[c]*Indiana University Indianapolis, Indianapolis, IN, USA*
[d]*Department of Otolaryngology-Head and Neck Surgery, Medical University of South Carolina, Charleston, SC, USA*
[e]*Department of Otolaryngology-Head and Neck Surgery, University of California, Los Angeles, Los Angeles, CA, USA*
[f]*Department of Otolaryngology-Head and Neck Surgery, University of Utah, Salt Lake City, UT, USA*
[g]*Department of Otolaryngology-Head and Neck Surgery, Oregon Health Sciences University, Portland, OR, USA*
[h]*Department of Otolaryngology-Head and Neck Surgery, Indiana University School of Medicine, Indianapolis, IN, USA*

## Abstract

Artificial intelligence (AI) has increasingly transformed medical data prognostics by enabling rapid and accurate analysis across imaging and pathology. However, investigation of machine learning predictions applied to prospectively collected, standardized, data from observational clinical intervention trials remains underexplored, despite its potential to reduce costs and improve patient outcomes. Chronic Rhinosinusitis (CRS), a persistent inflammatory disease of the paranasal sinuses that lasts more than three months, imposes a substantial burden on quality of life (QoL) and cost to society. Although many patients respond to medical therapy, others with refractory symptoms often pursue surgical intervention. Surgical decision-making in CRS is complex, as it must weigh known procedural risks against unclear individualized outcomes. In this study, we evaluated the utility of supervised machine learning models for predicting surgical benefit in CRS, using the main patient-reported outcome (Sino-Nasal Outcome Test-22, SNOT-22) as

the primary outcome. Our cohort of prospectively collected data from an observational intervention trial comprised patients who all underwent surgery; we specifically investigate whether models trained only on preoperative data could have identified those who might not have been recommended surgery prior to the procedure. Across multiple algorithms, including an ensemble approach, our best model achieved ∼85% classification accuracy, delivering accurate and interpretable predictions of surgical candidacy. Moreover, on a heldout set of 30 cases spanning mixed difficulty, our model achieved 80% accuracy, exceeding the average prediction accuracy of expert clinicians (75.6%), further demonstrating its potential to augment clinical decision-making and support personalized care in CRS management.

## 1. Introduction

Artificial Intelligence (AI) for healthcare is emerging as one of the most dynamic areas of research and development worldwide. The foundational ideas date to the 1950s, when the Turing test was proposed as an operational benchmark for machine intelligence. Subsequent methodological and hardware advances over the years have enabled translational applications across health sciences [1]. Today, AI and machine learning (ML) power a broad spectrum of products and services: from business analytics and robotics to voice-interactive assistants such as Siri, Alexa, and Google Assistant [2]. Healthcare, in particular, has seen rapid growth in AI investment and adoption since 2016, reflecting the potential of technology in this industry to improve outcomes and reduce costs [1, 3, 4, 5, 6].

There are many uses for ML/AI in medical applications for clinical cases beyond basic operations. For clinicians, high-quality treatment recommendations depend on deep domain expertise and experience in interpreting complex heterogeneous information. As electronic health records (EHRs) have become the primary medium of documentation, clinicians spend substantial time on data entry, yet the resulting information, collected in part to enable billing and administrative workflows and in part to support clinical care and longitudinal analysis, remains difficult to synthesize during brief patient encounters. In addition, EHR data span both structured elements (e.g.,

demographics, diagnoses, medications, laboratory values, and standardized patient-reported outcomes) and unstructured elements (e.g., free-text clinical notes), each offering complementary signals whose utility depends on the study context and modeling objective. In this work, we focus on structured preoperative clinical variables and standardized patient-reported outcomes, which are readily comparable across patients and amenable to reproducible model development. For treatment recommendations, risk stratification, and outcome forecasting in personalized medicine, ML models can help shoulder this cognitive load by learning patterns from large datasets and producing fast, reproducible risk estimates or recommendations [2]. In one possible collaborative paradigm, AI augments, not replaces, clinical judgment: physicians and patients focus their limited time on interpreting model-informed options, clarifying trade-offs, and engaging in shared decision-making. Although AI has transformed several diagnostic domains (e.g., imaging and digital pathology) [7], the next sequential phase of use in the treatment domain remains comparatively underexplored.

Introducing AI-based data analytics to inform treatment selection enables physicians to synthesize data and findings learned from large, heterogeneous clinical datasets rapidly at the point-of-care. Objective, data-driven modeling can strengthen patient trust, help avoid unnecessary procedures, and support high-value, cost-effective care [8, 9, 10]. Neither AI nor the most experienced clinician can achieve 100% accuracy; however, when used collaboratively, ML predictions can complement human expertise and improve decision quality, standardizing subjective data interpretation, potentially elevating less experienced non-specialty physicians [3, 11]. In this study, we report the results of supervised ML modeling in one of the most prevalent and costly health problems in the United States, Chronic Rhinosinusitis (CRS), which affects a substantial fraction of adults and is a leading driver of outpatient visits, medication use, and overall healthcare utilization [12, 13]. When symptoms remain refractory to appropriate medical therapy, endoscopic sinus surgery is commonly pursued with the expectation of meaningful improvement in sinonasal quality of life; however, responses are heterogeneous and a non-trivial subset of patients experience limited benefit despite procedural risk and cost. We evaluate supervised ML prediction models for supporting *surgical recommendation* in CRS using SNOT-22 as the reference outcome. Importantly, our prospectively collected cohort comprises patients who all *underwent surgery*; we specifically ask whether models trained solely on *preoperative* data could have *preemptively* identified those who might not have

been recommended surgery (limited SNOT-22 benefit = class 0), despite ultimately receiving it. We formulate this as a counterfactual classification task in which class 1 indicates "surgery recommended as per the outcome of the prediction model" and class 0 indicates "surgery not recommended as per the outcome of the prediction model," with labels derived from 6-month SNOT-22 change relative to a minimal clinically important difference threshold [14, 15, 16].

We benchmark a suite of models, including logistic regression, support vector machines, random forests, Naïve Bayes, and neural networks, as well as ensemble strategies. Among these, a compact *three-layer DNN (one hidden layer)* and a *majority-voting ensemble* performed best, achieving $\sim$85% classification accuracy. The main contributions of this work are:

- We frame CRS surgical selection as a preoperative surgical benefit prediction task, estimating the likelihood of achieving clinically meaningful improvement following ESS, training solely on *pre-operative* variables.

- We implement and compare multiple ML models and ensembles, demonstrating that AI is feasible for CRS surgery recommendation. To our knowledge, this is the first large-scale, multi-center study of its kind.

- We conduct experiments to benchmark the performance of AI based algorithms to those of 6 human experts. Our results show that the proposed approach outperforms the experts on average in CRS surgical outcome prediction.

## 2. Chronic Rhinosinusitis

*2.1. Economics on CRS*

Chronic rhinosinusitis (CRS) is a persistent inflammatory disorder of the nose and paranasal sinuses, defined clinically by a constellation of local symptoms lasting $\geq$ 3 months (Fig. 1) [13]. In Western adult populations, its prevalence is commonly estimated at roughly 10–15% [13]. Although acute infections may precipitate or accompany disease flares, the hallmarks of CRS are *chronic mucosal inflammation* with tissue remodeling and associated disturbances of mucus clearance and epithelial function [18]. Symptom burden is substantial: patients frequently report facial pain/pressure and headache, nasal obstruction and congestion, rhinorrhea/postnasal drainage, and loss
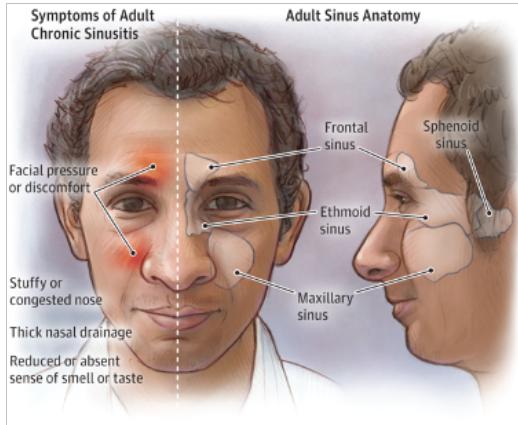
Figure 1: Classic local manifestations of Chronic Rhinosinusitis (CRS). In addition, systemic such as fatigue and depression, and financial burden of direct and indirect costs, together result in a large overall disease burden (Figure reused with permission from [17]).

of smell/taste, while systemic or non-nasal manifestations (e.g., ear pain, dental pain, cough, sleep disruption, fatigue, and mood symptoms) further compound individual quality-of-life (QoL) impairment. Beyond sinonasal symptoms, CRS is associated with broad health-related quality-of-life declines, exacerbation of comorbid atopic and respiratory diseases, work and school through absenteeism and presenteeism, and commonly involves repeated healthcare encounters and courses of antibiotics and systemic corticosteroids with attendant adverse effects and costs [12, 13, 19, 20, 21, 22, 23]. Overall QoL decrements in CRS are comparable to those seen with major chronic conditions such as congestive heart failure, chronic obstructive pulmonary disease, end-stage renal disease on dialysis, and chronic low back pain [24]. The societal impact is likewise considerable, with costs accruing from direct medical care and productivity losses; per-patient annual economic burden in the United States has been estimated at approximately $10,077 [12].

## 2.2. CRS Disease Management

Initial treatment consists of medical treatments, including topical saline lavage, intranasal or oral corticosteroids, and antibiotic therapy. If these or further medical therapies result in significant symptom improvement, patients may continue such treatments for the long-term to successfully manage their disease. However, many patients do not gain sufficient benefit from

medical treatments and are offered the option of Endoscopic Sinus Surgery (ESS) [14].

## 2.3. Significance of CRS Surgical Outcome Prediction

Analyses of the National Ambulatory Medical Care Survey (NAMCS) from the late 1980s and early 1990s reported rising office visits for sinusitis and identified sinusitis as one of the leading indications for antibiotic prescribing in U.S. outpatient practice [25]. More recently, a national claims analysis showed that rhinosinusitis accounts for *more outpatient antibiotic prescriptions than any other diagnosis*, underscoring its outsized contribution to community antibiotic use [26]. Chronic rhinosinusitis (CRS) is also a major economic burden: contemporary estimates place direct medical costs in the United States at roughly \$10–\$13 billion annually, with an additional $\sim$ \$20 billion in indirect costs from productivity loss [12, 27].

Given these statistics, it is paramount to achieve rapid and durable disease control in this chronic disease by quickly traversing the clinical management algorithm. In CRS, patient-reported quality of life (QoL) is a primary outcome and decision anchor, with instruments such as the SNOT-22 recommended for tracking severity and response [16, 13, 28]. According to current estimates, about 300,000 ESS operations for CRS have been carried out annually in the U.S, with 30–40% of patients showing unsatisfactory results in terms of achieving clinically significant SNOT-22 benefits [14, 15]. According to another recent study, 20-30% of CRS patients who underwent endoscopic sinus surgery (ESS) reported that their post-treatment symptom reduction fell short of their expectations. Patient satisfaction with CRS intervention was strongly associated with the extent to which outcome expectations were met. [29].

Most patients with medically refractory CRS do improve with surgery, but a substantial minority, on the order of one in three, derive limited benefit when SNOT-22 is the primary outcome of interest [14, 15]. This variability at the individual level motivates our study: if clinicians and patients could access *pre-operative*, patient-specific estimates of surgical benefit, they could better weigh procedural and anesthesia risks against likely QoL gains, engage in more informed shared decision-making, and set realistic expectations [11]. Such individualized risk-benefit estimation is a central promise of personalized medicine and a key motivation for the rapid growth of structured data collection in healthcare over the past decade. Leveraging current AI/ML

6

methods, we therefore investigate whether models trained on routinely available clinical data can estimate surgical outcomes to support counseling for elective ESS in CRS.

## 3. Related & Prior Work

### 3.1. Clinical guidance for CRS decision-making

Contemporary care pathways for chronic rhinosinusitis (CRS) emphasize accurate phenotyping, optimization of medical therapy, and patient-centered shared decision-making about elective endoscopic sinus surgery (ESS) when symptoms remain refractory [13, 30]. Patient-reported outcomes (PROs)–in particular the SNOT-22–are recommended to quantify baseline burden and to track response to therapy [16, 13]. Several longitudinal studies have demonstrated efficacy of ESS in this situation, and identified risk factors for treatment success/failure. But, the medical and sociodemographic variables are complex and interactive, and to date, no risk calculator exists in practice. Meta-analyses consistently show large average post-ESS improvements in SNOT-22, yet with wide interpatient variability [14, 15]. Current guidelines provide indication frameworks, but they stop short of offering individualized, pre-operative predictions of benefit using routinely available clinical variables [30]. This creates a practical gap for clinicians and patients seeking personalized risk–benefit estimates to guide shared decision-making.

### 3.2. Predictive models in rhinology and sinus surgery

Predictive modeling in rhinology remains comparatively sparse and heterogeneous. Prior efforts include regression-based and machine-learning (ML) approaches that associate baseline symptom burden and selected clinical variables with post-operative SNOT-22 trajectories or satisfaction; for example, Kang et al. used SNOT-22 to model outcomes after septoplasty [31], while large meta-analyses quantified average gains but were not designed as individualized risk calculators [14, 15, 29]. Much of the existing literature centers on population-level effect sizes, disease control categories, or biomarker/imaging correlates rather than point-of-care prediction from routine pre-operative clinical records [13]. Consequently, there is no widely adopted, validated pre-operative tool that estimates a CRS patient's probability of achieving a clinically meaningful SNOT-22 improvement after ESS using standard clinical data elements. A recent single-center study [32] trained machine-learning models to predict attainment of SNOT-22 MCID after primary ESS using 59

7

preoperative predictors collected via an institutional survey platform, and reported strong discriminative performance for a stacking-style Ensemble model (AUC up to 0.89; accuracy 87.8%) on an 80:20 split of 242 patients (test $n = 48$). While promising, the study's clinical generalizability remains uncertain because the cohort is drawn from a single academic sinus center and evaluation relies on a relatively small held-out test set. In addition, postoperative SNOT-22 outcomes were permitted as early as $\geq 2$ months (when stability evidence is weaker than at $\geq 6$ months), and the TabNet component required assigning specific numeric values to represent missing data, both of which may introduce additional variance and warrant careful sensitivity analysis. Overall, more rigorous validation (e.g., multicenter external validation and stronger deployment-oriented evaluation) is still needed before such models can be relied upon for point-of-care surgical counseling.

### 3.3. Machine learning for surgical decision support

Across surgery more broadly, interpretable ML is increasingly explored to augment case selection, perioperative risk stratification, and shared decision-making [11, 7, 3]. Also transparency in the process from data gathering, processing, assessment, experiments, reporting etc. is now expected, using TRIPOD+AI guidelines as we do here [33] or similar alternative (CHAI [34]). A paper in this arena should follow at least one of these and mention which one. Best practices emphasize calibration, transparency of feature influence, and evaluation of net clinical utility versus treat-all or treat-none heuristics [11]. However, translation to CRS has lagged: most published clinical ML algorithms either rely on imaging-heavy pipelines or small, single-center cohorts, and they frequently omit core translational checks such as probability calibration and decision-curve (net-benefit) analysis—steps that are essential for clinical deployment [31, 11, 35, 36, 37, 38].

### 3.4. Gap addressed in this work

In sum, while guidelines endorse PRO-driven management and meta-analyses demonstrate average post-ESS improvements, there remains a clear gap: a systematic ML application to *routinely accessible* CRS clinical data that (i) frames the task as a pre-operative recommendation problem anchored to the SNOT-22 minimal clinically important difference, (ii) compares multiple algorithms and ensembles, (iii) prioritizes interpretability and visualization, and (iv) quantifies potential clinical utility compared to experts. Our study addresses these needs by building and evaluating such models on a

multicenter CRS surgery cohort, and we report our approach and evaluation in a systematic fashion according to TRIPOD+AI recommendations.

## 4. Dataset

For predicting surgical outcomes (i.e., the potential effectiveness of ESS in CRS patients considering surgery), training the model is a major initial task. Currently, the Sino-Nasal Outcome Test-22 (SNOT-22) is the most widely accepted sinus QOL tool for analyzing CRS patient-reported outcomes. It is an extensively validated set of 22 standard questions for understanding CRS symptom severity [14, 15, 28, 29]. Here, the patient rates each question on a Likert-type scale in the range of 0 to 5, resulting in a total score range between 0 to 110. Higher scores depict the worst QoL [15].

To understand outcomes from medical and surgical intervention in CRS, an NIH-funded multicenter CRS outcomes study was designed and undertaken (clinicaltrials.gov NCT01332136). This study recorded detailed demographic and clinical metadata in a longitudinal fashion, along with patient-reported outcomes, including (1) a sinus-specific evaluation (SNOT-22), and (2) a general health measure (SF6d-derived health utility value, HUV) [39]. A standardized collection of ∼30 clinical and demographic common data elements available to practicing physicians was performed throughout this study and used as the set of predictor variables.

We began with a development cohort of 791 patients (50 attributes; typical missingness for some clinical fields, referred to as the 2R01 dataset for the rest of the paper) and a second, independent cohort of 355 patients from a related multicenter CRS outcomes study with comparable clinician-recorded metadata (henceforth, referred to as the 3R01 dataset). In this paper, we restrict analysis to patients who actually underwent sinus surgery, because our objective is to build a model that flags cases that may be recommended to avoid surgery (i.e., would not achieve sufficient postoperative benefit). Operationally, this means we considered only the "surgery" cases, excluded patients managed medically (no surgery). So, we concatenated the two cohorts, and retained only those with sinus surgery as the treatment. After these filters, the combined dataset comprised 524 surgical patients. For this multicenter, observational CRS surgical outcomes study with six different tertiary rhinology practice sites, institutional review board approval was gained for each site, and de-identified data were shared for this investigation (COMIRB #19-2085).

9

The primary goal in this analysis is to be able to predict whether a surgery would result in the desired postoperative outcome (yes/no). To obtain that goal, we need to predict the change in total SNOT-22 from baseline to six months post-ESS, the point at which the treatment outcome is statistically considered stable and durable [31]. The minimal clinically important difference (MCID) is a metric used to estimate the minimal amount of change in the SNOT-22 measurement that is noticeable clinically and is used to describe the clinical importance that results from a CRS intervention. If the SNOT-22 score improvement is greater than or equal to an established MCID (8.9 points), it indicates that the intervention (ESS) resulted in a noticeable gain in sinus-related QoL [31]. If the SNOT-22 score change was less than 8.9 points, then surgery did not result in noticeable benefit and perhaps continued medical therapy could have been continued unless there were other reasons to pursue surgical intervention. This categorization is used as a major binary output classification for the model training in the current paper.

## 5. Methodology

### 5.1. Supervised learning paradigm

Machine learning (ML) systems are commonly organized into three broad paradigms: *supervised* learning, *unsupervised* learning, and *reinforcement* learning [40, 41, 42]. In supervised learning—the focus of this study—a model is learned from labeled examples $(x_i, y_i)$ to approximate an unknown mapping $f : \mathcal{X} \to \mathcal{Y}$, such that

$$y \approx f(x),$$

where $x$ denotes an input feature vector (e.g., demographics, comorbidities, baseline SNOT-22 items) and $y$ is the associated target label (here, a binary indicator of surgical recommendation). The objective is not merely to fit the training data but to *generalize* to unseen cases drawn from the same data-generating process [41, 40].

Model development typically involves (i) splitting data into training/test sets or using cross-validation for hyperparameter tuning, (ii) fitting candidate algorithms, and (iii) assessing performance on the held-out test data with appropriate metrics [41]. For binary clinical prediction tasks, accuracy alone can be misleading, particularly under class imbalance; complementary metrics such as precision, recall, F1-score, area under the ROC curve (AUROC), and confusion matrix are recommended [43, 37].
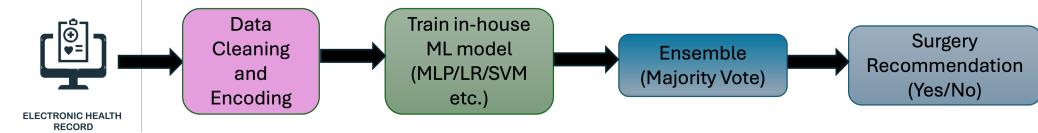
Figure 2: Schematic pipeline of the proposed decision-support approach. Structured fields from the electronic health record (EHR) undergo data cleaning and encoding, after which multiple in-house classifiers (e.g., MLP, logistic regression, SVM) are trained. Model outputs are combined via a majority-vote ensemble to generate a binary surgery recommendation (Yes/No).

In this work, we framed CRS surgical selection as a supervised, preoperative classification problem and compared multiple algorithms including regularized logistic regression, SVMs, random forests, gradient-boosted trees, Naïve Bayes, a compact three-layer DNN, and a majority-vote ensemble, using stratified validation and a held-out test set. We report discrimination (e.g., AUROC), and error profiles (e.g., class-0 recall and macro-F1), consistent with best practices for clinical prediction modeling [37, 35, 11]. Regarding patient and public involvement, neither was involved in the study design or any feedback collection.

### 5.2. Machine Learning Approaches Utilized

We implement the binary classification models using `scikit-learn` [44] and `XGBoost` [45] libraries. In Stage 1, we train and benchmark individual classifiers (logistic regression, Naïve Bayes, support vector machine, random forest, multi-layer perceptron, and XGBoost). In Stage 2, we form an ensemble by aggregating model outputs to improve predictive performance. Unless otherwise noted, hyperparameters were tuned via cross-validation, and performance was assessed on a held-out test set. The whole pipeline is schematically depicted in Fig. 2.

### 5.2.1. Logistic Regression

Logistic regression models the conditional probability of the positive class via the sigmoid of an affine function [46]:

$$p(y = 1 \mid \mathbf{x}) = \sigma(z), \qquad z = \mathbf{w}^\top \mathbf{x} + b, \qquad \sigma(z) = \frac{1}{1 + \exp(-z)}. \quad (1)$$

11

Parameters $(\mathbf{w}, b)$ are estimated by minimizing the regularized negative log-likelihood:

$$\min_{\mathbf{w}, b} -\sum_{i=1}^{n} \left[ y_i \log \sigma(z_i) + (1 - y_i) \log(1 - \sigma(z_i)) \right] + \lambda \mathcal{R}(\mathbf{w}), \qquad (2)$$

where $\mathcal{R}(\mathbf{w})$ is typically $\ell_2$ (ridge) or $\ell_1$ (lasso) regularization.

*5.2.2. Naïve Bayes*

Naïve Bayes applies Bayes' rule with a conditional independence assumption over features given the class [47]:

$$p(y \mid \mathbf{x}) \propto p(y) \prod_{j=1}^{d} p(x_j \mid y), \qquad (3)$$

with the decision rule $\hat{y} = \arg\max_{y \in \{0,1\}} \left\{ \log p(y) + \sum_{j=1}^{d} \log p(x_j \mid y) \right\}$.

*5.3. Random Forest*

Random forests are ensembles of decision trees trained on bootstrap samples with feature subsampling at each split [48]. It creates decision trees from several observations. For classification, the forest prediction is the majority vote (or averaged class posterior) over $T$ trees $\{h_t\}_{t=1}^{T}$:

$$\hat{y} = \arg\max_{k \in \{0,1\}} \frac{1}{T} \sum_{t=1}^{T} \mathbb{I}\{h_t(\mathbf{x}) = k\}, \qquad (4)$$

and uses the average for regression problems. In decision trees, each step involves a stingy selection of the optimal split point from the training dataset. By building several trees using various sampled datasets from the training datasets can decrease the high variance. This process is known as bootstrap aggregation, or bagging. Bagging reduces variance, and random feature selection decorrelates trees, improving generalization.

*5.4. Support Vector Machine (SVM)*

This method is frequently applied for classification problems. SVM is effective in high-dimensional spaces. This algorithm projects each data sample as a point in an n-dimensional space, where n is the number of features.

Next, classification is carried out by identifying the hyper-plane that effectively distinguishes the two classes. The selection of hyperplane is based on maximizing the distance between the closest data points of each class. SVM also uses different kernels, which can convert the low dimensional space into high dimensional spaces when the problem is non-linear. For linearly separable data, the hard-margin SVM maximizes the geometric margin [49]. In the soft-margin setting:

$$\min_{\mathbf{w},b,\boldsymbol{\xi}} \quad \frac{1}{2}\|\mathbf{w}\|_2^2 + C\sum_{i=1}^{n}\xi_i \tag{5}$$

$$\text{s.t.} \quad y_i(\mathbf{w}^\top\phi(\mathbf{x}_i)+b) \geq 1-\xi_i, \quad \xi_i \geq 0, \; i=1,\ldots,n, \tag{6}$$

where $\phi(\cdot)$ is an implicit feature map induced by kernel $K(\mathbf{x},\mathbf{x}') = \langle\phi(\mathbf{x}),\phi(\mathbf{x}')\rangle$. The decision function is

$$f(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^{n}\alpha_i y_i K(\mathbf{x}_i,\mathbf{x})+b\right), \tag{7}$$

with dual coefficients $\alpha_i$ found by quadratic programming.

*5.5. Multi-Layer Perceptron (MLP)*

Usually, Artificial Neural Networks are referred to as neural networks or Multi-Layer Perceptrons since it is the most widely used architecture in the field of neural networks. A perceptron incorporates a single neuron model, which is a node in larger neural networks. The layered structure of the framework of the network gives the power of prediction for neural networks. The layered architecture consists of the input layer, hidden layers, and an output layer. An $L$-layer MLP composes affine maps and nonlinearities [50]:

$$\begin{aligned} \mathbf{h}^{(0)} &= \mathbf{x}, \quad \mathbf{h}^{(\ell)} = \sigma^{(\ell)}\big(\mathbf{W}^{(\ell)}\mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)}\big), \; \ell = 1,\ldots,L-1, \\ \hat{p} &= \text{sigmoid}\big(\mathbf{w}^\top\mathbf{h}^{(L-1)}+b\big). \end{aligned} \tag{8}$$

trained by minimizing binary cross-entropy (optionally with weight decay) via stochastic gradient descent, with the loss function being:

$$\mathcal{L} = -\frac{1}{n}\sum_{i=1}^{n}[y_i\log\hat{p}_i + (1-y_i)\log(1-\hat{p}_i)] + \lambda\sum_{\ell}\|\mathbf{W}^{(\ell)}\|_2^2. \tag{9}$$

## 5.6. Extreme Gradient Boosting (XGBoost)

XGBoost builds an additive ensemble of regression trees to minimize a regularized objective using second-order (Newton) boosting [45, 51]. Let $F_M(\mathbf{x}) = \sum_{m=1}^{M} f_m(\mathbf{x})$ be the model, where each $f_m$ is a tree. The objective is

$$\mathcal{L}^{(M)} = \sum_{i=1}^{n} \ell\Big(y_i, \hat{y}_i^{(M-1)} + f_M(\mathbf{x}_i)\Big) \ + \ \Omega(f_M), \qquad \Omega(f) = \gamma T + \frac{\lambda}{2}\|\mathbf{w}\|_2^2, \tag{10}$$

where $T$ is the number of leaves, $\mathbf{w}$ are leaf scores, and $\ell$ is the logistic loss for binary classification. Using a second-order Taylor expansion around $\hat{y}_i^{(M-1)}$, the split gain for a candidate partition $(L, R)$ with gradients $g_i = \partial_{\hat{y}}\ell$ and Hessians $h_i = \partial_{\hat{y}}^2\ell$ is

$$\text{Gain} = \frac{1}{2}\left(\frac{\left(\sum_{i \in L} g_i\right)^2}{\sum_{i \in L} h_i + \lambda} + \frac{\left(\sum_{i \in R} g_i\right)^2}{\sum_{i \in R} h_i + \lambda} - \frac{\left(\sum_{i \in L \cup R} g_i\right)^2}{\sum_{i \in L \cup R} h_i + \lambda}\right) - \gamma. \tag{11}$$

Trees are added greedily to maximize (11) until convergence or early-stopping criteria are met.

*Ensemble aggregation..* Given $M$ trained models producing class posteriors $p_m(y = k \mid \mathbf{x})$, we use soft voting with nonnegative weights $\{w_m\}_{m=1}^{M}$ (normalized to sum to 1):

$$\hat{y} = \arg\max_{k \in \{0,1\}} \sum_{m=1}^{M} w_m\, p_m(y = k \mid \mathbf{x}). \tag{12}$$

We also experimented with stacking (logistic meta-learner) using out-of-fold predictions as meta-features [52]. In particular, the majority voting ensemble process during inference is illustrated schematically in Fig. 3. As shown in Fig. 3, the pipeline ingests structured clinical variables, performs data cleaning and encoding, and feeds the resulting features in parallel to multiple base learners (`Random Forest`, `Logistic Regression`, `SVM`, `MLP` etc.). Their predictions are aggregated via a voting scheme to yield a binary surgery recommendation (Yes/No). Only preoperative variables are used as features for training to prevent information leakage; the ensemble reduces variance and mitigates model-specific biases, providing more stable recommendations than any single classifier.
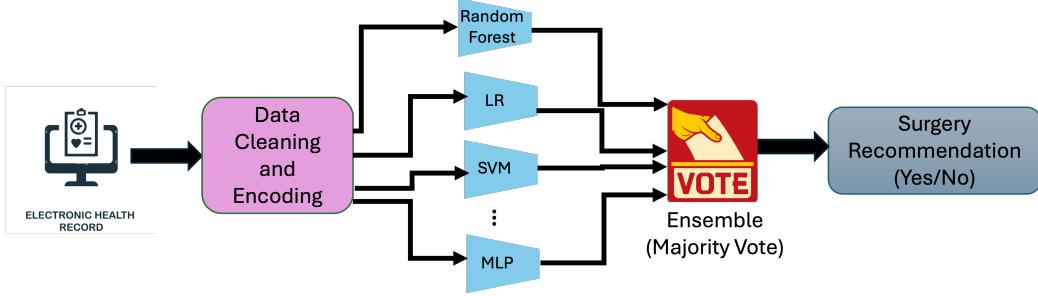
Figure 3: Ensemble decision-support pipeline. Structured EHR data are cleaned and encoded, then passed in parallel to multiple base learners (`Random Forest`, `Logistic Regression`, `SVM`, `MLP` etc.). Their predictions are aggregated via a voting scheme to produce a binary surgery recommendation (Yes/No). Only preoperative variables are used for inference.

## 6. Dataset Preparation

### 6.1. Cleaning and Encoding

Dataset cleaning is an important step for every machine learning model. This is of particular relevance for clinical data, as there are often challenges in the systematic measurement and recording of data that result from varied human factors. A messy dataset might negatively impact the prediction model, so it is necessary to examine and understand any flaws and address them wherever possible before training a model alongside reporting this process transparently. Here, this was accomplished in two parts. First, the data were cleaned and then preprocessed (encoding, then categorical to numerical conversion) to obtain a set of features suitable for ML training. Note, this process was repeated for both the 2R01 and 3R01 datasets.

For the 2R01 dataset, initially, there were a total of 791 rows of patient data with 50 attributes as columns. First, we chose the subset where treatment was sinus surgery, reducing the dataset to 604 rows. There were missing data and unwanted attributes present which might not be useful for accurate surgery outcome prediction. So, in the cleaning process, all the unwanted columns were removed from this dataset under a physician's guidance. Then, all the blank spaces (i.e., the null values) in all the columns of the patient were removed. Then the total null value counts were found as shown in Table 1.

To establish the target outcome, the difference between Baseline SNOT-22 and 6-month postoperative SNOT-22 scores was calculated for each in-

Table 1: Feature-wise missing (null) values in the 2R01 cohort restricted to surgery-treated patients.

| Feature | # Nulls |
|---|---|
| Age | 1 |
| Race | 1 |
| Education | 10 |
| Household Income | 7 |
| Smoker | 4 |
| Alcohol | 5 |
| Diabetes | 1 |
| Baseline CT Score | 2 |
| Baseline Endoscopy Score | 3 |
| SNOT-22 Baseline Score | 1 |

dividual. After removing the 217 patient samples where patients did not come back for followup after 6 months (as such, we do not have a label for them), baseline Health Utility score columns were also removed since the study concentrated on the SNOT-22 score. Now to find the difference between the baseline SNOT-22 score and the SNOT-22 score after 6 months of operation, we changed their datatype into integers since the values in those columns were object types. Then, the differences between those two columns were calculated and assigned as the target column. Additionally, the 6 month SNOT-22 score column was removed as it is a post-operative attribute. Also, the rest of the null value rows (Table 1) were removed as they were very few in number. Then the dimensions of the dataset became (371,31). Overall, the missingness of the fields in the dataset was 4.7%, which is quite low. We also experimented with data imputation for the missing values, however, as these cases were very few (3% only), it did not affect the results compared to excluding the subject (sample) altogether. Similar data cleaning was performed on the 3R01 dataset, which initially had 354 rows of patient data with 39 attributes as columns. Zooming in on just the surgery treatment cases resulted in 266 rows. The null values in specific columns are shown in Table 2. Similar operation was carried out to create the targets, which resulted in a cleaned dataset with dimensions (153,31). The attributes for this 3R01 dataset were organized to overlap with those in the 2R01 dataset. As in the 2R01 dataset case, the target column for the 3R01 dataset was also

Table 2: Feature-wise missing (null) values in the 3R01 cohort restricted to surgery-treated patients.

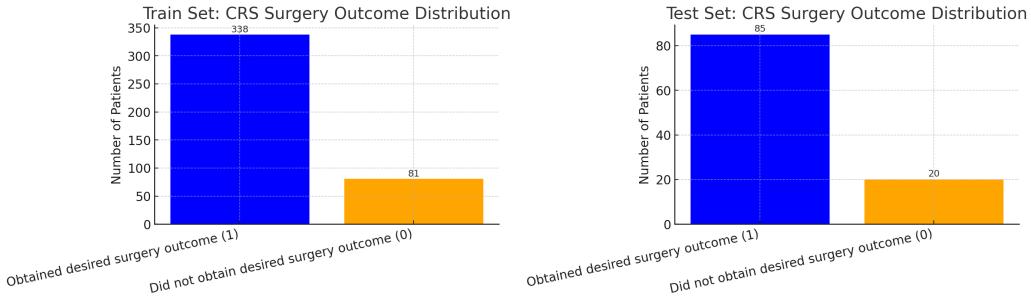| Feature | # Nulls |
|---|---|
| Race | 2 |
| Education | 3 |
| Household Income | 1 |
| Obstructive Sleep Apnea (OSA) History | 3 |
| Alcohol | 3 |
| Baseline Endoscopy Score | 3 |
| SNOT-22 Baseline Score | 3 |



Figure 4: Class distribution in the train (left) and test (right) sets. Bars show the number of CRS patients who obtained the desired surgery outcome (1) versus those who did not obtain the desired outcome (0). The splits retain the overall class imbalance (train: 338 vs. 81; test: 85 vs. 20), which is considered in model training and evaluation.

changed from continuous values to binary outcomes based on the established SNOT-22 8.9 point MCID.

Next, the categorical variables were converted to numeric codes using deterministic, clinically sensible mappings. This encoding is necessary because most machine-learning algorithms operate on numerical tensors and require consistent, reproducible representations of categorical information to learn stable decision boundaries. Specifically, demographic/socioeconomic fields (e.g., `Sex`, `Race`, `Ethnicity`, `Education`, `Household Income`, `Insurance`) and clinical history variables (e.g., `Previous Surgery`, `Positive allergy testing` , `Smoker`, `Alcohol`, `Chronic Obstructive Pulmonary Disease (COPD)`, `Aspirin Intolerance`, `Obstructive Sleep Apnea (OSA) History`, `Gastroesophageal Reflux Disease (GERD)`, `Allergic Fungal Sinusitis`

(AFS), `Asthma`, `Diabetes`, `Fibromyalgia`) were recoded from text labels into integer categories via explicit dictionaries (for example, *Female*→0, *Male*→1 for `Sex`; *Employer provided*→0, *Medicare*→1, *Private*→2, *Canadian Medicare*→3, *Medicaid*→4 for `Insurance`). Continuous baseline measures, such as the `SNOT-22 Baseline Score`, `Baseline CT Score`, and `Baseline Endoscopy Score`, were retained as numeric features. Where applicable, missing values were handled using simple domain-consistent rules (e.g., mapping textual "None" to a sentinel category for certain comorbidities); otherwise they were left for model-specific handling. After the encoding, we verified type consistency, and exported the encoded dataset for downstream train/test splitting and model fitting. This encoding scheme enables reproducible training across classifiers and avoids leakage from post-treatment or follow-up fields.

Next, we partitioned the dataset into an **80:20 train–test split** using stratified sampling on the binary outcome to preserve class prevalence. **Stratified sampling** here means that we performed the split *within each outcome class* (class 1 and class 0) and then combined the selected patients, yielding similar class proportions in both subsets. The training set comprised 338 patients with the desired postoperative outcome (class 1) and 81 without (class 0); the test set comprised 85 class 1 and 20 class 0 patients. This yields closely matched prevalence across splits (train: 80.6% vs. 19.3%; test: 81.0% vs. 19.0%). As shown in Fig. 4, the side-by-side bar plots depict these counts, with blue indicating patients who obtained the desired outcome (1) and orange indicating those who did not (0). Stratification enables fair comparison between training and held-out performance.

*6.2. Feature Correlation Analysis*

To explore linear structure and potential dimensionality reduction, we applied Principal Component Analysis (PCA) to the standardized feature matrix (zero mean, unit variance; no outcome or post-treatment variables included). The explained-variance plot (top left in Fig. 5) shows that the first few components account for a disproportionate share of variance, followed by a long tail of smaller contributions. The cumulative explained-variance curve (top right in Fig. 5) indicates that retaining a moderate number of principal components captures the vast majority of variance (the 95% reference line is shown), which is useful for building compact visualizations or low-dimensional baselines while avoiding information leakage. Finally, the PC1–PC2 scatter (colored by class) reveals substantial overlap between

groups in the first two linear components (bottom right in Fig. 5), suggesting that class separation, if present, likely resides in higher-order directions and/or is non-linear. This motivates the use of flexible learners (e.g., tree ensembles, kernels, or neural networks) and careful regularization rather than relying solely on linear boundaries in a low-dimensional subspace. Note, PCA here is used for exploratory analysis and visualization; model training uses the original engineered/encoded features with appropriate preprocessing. Furthermore, we examined pairwise linear associations among all preprocessed clinical predictors using the Pearson correlation coefficient after encoding categorical variables and standardizing continuous features, as shown in bottom left plot of Fig. 5. The heatmap (blue = negative, red = positive) visualizes the full correlation matrix, with the diagonal fixed at 1 and off-diagonal cells indicating the strength and direction of association between feature pairs. This screening step helps identify redundancies and potential collinearity prior to modeling and interpretation. While tree-based learners (e.g., Random Forest, XGBoost) are relatively robust to correlated inputs, we still monitor high absolute correlations to avoid unstable coefficients in linear models (e.g., logistic regression) and to interpret variable importance more cautiously. No outcome or post-treatment variables were included in the matrix to prevent leakage.

## 7. Results and Analysis

*First Level Results*

In the first level, the training split of the combined dataset was used to individually train the following classifiers: Logistic Regression, Multi-Layer Perceptron (MLP), Random Forest (RF), Support Vector Machine (SVM,) and Naïve Bayes. Then, the testing split (20% of data) was used for evaluating the trained model. Note, once the models were trained, they were not updated based on new data. The results are presented below:

*7.0.1. Logistic Regression*

As shown in the **top row** of Fig. 6, logistic regression attains overall accuracy **0.85**. Class 1 (desired outcome) performance is strong (precision **0.86**, recall **0.98**, F1 **0.91**), while class 0 remains challenging (precision 0.75, recall 0.30, F1 0.43). The confusion matrix $[6, 14; 2, 83]$ indicates a tendency to predict class 1, prioritizing sensitivity over specificity.

### 7.0.2. Support Vector Machine (SVM)

As shown in the **second row** of Fig. 6, SVM achieves accuracy **0.79**. It performs well on class 1 (precision **0.82**, recall **0.95**, F1 **0.88**) but exhibits low recall for class 0 (0.10; F1 0.15). The confusion matrix $[2, 18; 4, 81]$ reflects many class $0 \rightarrow 1$ errors, consistent with the class imbalance.

### 7.0.3. Naïve Bayes

As shown in the **third row** of Fig. 6, Naïve Bayes underperforms with accuracy **0.30**. It perfectly recalls class 0 (precision 0.21, recall **1.00**, F1 0.35) but misses most positives (class 1 recall 0.13, F1 0.23). The confusion matrix $[20, 0; 74, 11]$ suggests the conditional independence assumptions are poorly matched to this dataset.

### 7.0.4. Multi-Layer Perceptron (MLP)

As shown in the **fourth row** of Fig. 6, the MLP reaches accuracy **0.85**. It balances precision and recall for class 1 (precision **0.88**, recall **0.94**, F1 **0.91**) and improves class 0 detection relative to linear/kernal baselines (recall 0.45, F1 0.53). The confusion matrix $[9, 11; 5, 80]$ indicates better control of false positives than SVM.

### 7.0.5. Random Forest

As shown in the **bottom row** of Fig. 6, Random Forest attains accuracy **0.82**. Positive-class performance is strong (precision **0.84**, recall **0.95**, F1 **0.90**), while class 0 recall remains modest (0.25; F1 0.34). The confusion matrix $[5, 15; 4, 81]$ highlights a minority-class miss pattern similar to other models, though tree ensembling yields robust overall discrimination.

### 7.0.6. Hyperparameter tuning of the MLP

All candidate models were hyperparameter optimized via stratified cross-validation on the training split. Because the MLP achieved the best trade-off for our primary goal (detecting class 0), we report its tuning here. We first varied network depth and observed that architectures with more than one hidden layer consistently underperformed a single-hidden-layer model, particularly on the class 0 F1. Given the limited dataset size, this behavior is consistent with overfitting from excess capacity. We therefore fixed a *single hidden layer* and tuned its width. As shown in Fig. 7, overall accuracy and weighted F1 were relatively stable ($\approx 0.82$–$0.85$) across widths, whereas the *class 0 F1* varied more markedly, improving with larger widths

and peaking for wide layers (e.g., $\sim$ 400–480 units). The class 1 F1 remained high ($\approx$ 0.88–0.92) throughout. We selected a single-hidden-layer MLP with 400 neurons, combined with early stopping, class weighting, and $L_2$ regularization, to maximize class 0 detection while maintaining strong overall performance. Figure 8 summarizes the optimized MLP's performance on the held-out test set. Class-wise metrics (left) highlight the asymmetry between classes - high precision/recall for class 1 ($F_1 = 0.91$) and moderate but clinically useful performance for class 0 ($F_1 = 0.53$), while the macro/weighted averages (right) remain close to the overall accuracy (0.85), indicating stable behavior without excessive overfitting to the majority class.

### 7.0.7. Comparison among the Models

Because our clinical objective is to *prospectively flag patients unlikely to achieve the desired postoperative outcome* (class 0), we prioritize models with strong class 0 detection rather than overall accuracy alone. As shown in Fig. 9, the MLP offers the best trade-off for this goal: it attains the highest class 0 F1-score (**0.53**) with improved recall (**0.45**) relative to logistic regression (F1 0.43, recall 0.30), SVM (F1 0.15, recall 0.10), and random forest (F1 0.34, recall 0.25), while maintaining strong performance on class 1 (F1 0.91) and high overall accuracy (0.85). Although Naïve Bayes recalls all class 0 cases (recall 1.00), it severely underperforms on class 1 (recall 0.13) and overall accuracy (0.30), making it unsuitable for balanced clinical use. Accordingly, the MLP is the preferred model for downstream analyses and prospective decision support.

### 7.0.8. XGBoost

We trained gradient-boosted decision trees (XGBoost) and tuned hyperparameters via a stratified grid search on the training split. The optimal configuration was: {`colsample_bytree` = 1.0, `learning_rate` = 0.05, `max_depth` = 3, `n_estimators` = 200, `subsample` = 0.8}. As summarized in Fig. 10, the tuned model achieved overall accuracy $\approx$ 0.83 on the held-out test set with confusion matrix $\begin{bmatrix} 5 & 15 \\ 3 & 82 \end{bmatrix}$. Class-wise metrics show strong performance for class 1 (precision 0.85, recall 0.96, $F_1 = 0.90$), but more limited detection of class 0 (precision 0.62, recall 0.25, $F_1 = 0.36$). Because our clinical objective prioritizes identifying class 0 patients, the XGBoost model underperforms the MLP, which attains a higher class 0 $F_1$ (0.53) while maintaining comparable overall accuracy (0.85). We therefore retain

the MLP as the preferred model for downstream analysis and decision support.

*Second Level Results: The Ensemble Approach*

Combining complementary classifiers in an ensemble can improve generalization by averaging out model–specific errors and leveraging diversity among base learners, often outperforming any single model. Motivated by this, we implemented and evaluated *three* ensemble strategies, described in the next subsections, to assess whether aggregating heterogeneous learners yields more reliable predictions for our task.

### 7.0.9. Majority-voting ensemble

To assess whether combining complementary inductive biases could further stabilize predictions, we implemented a simple *hard-voting* ensemble using six independently tuned base learners: Logistic Regression, SVM, Random Forest, Naïve Bayes, MLP, and XGBoost (see individual results above). For each test instance $x_i$, the ensemble aggregates the class labels $\hat{y}_i^{(m)} \in \{0, 1\}$ from model $m = 1, \ldots, 6$ via majority vote

$$\hat{y}_i^{(\text{ens})} \;=\; \arg\max_{c \in \{0,1\}} \sum_{m=1}^{6} \mathbf{1}\Big(\hat{y}_i^{(m)} = c\Big),$$

with a deterministic tie–break in favor of the **MLP** prediction (consistent with our objective of maximizing class 0 detection and with its strong standalone performance).

As shown in Fig. 11, the majority-vote ensemble attains an overall accuracy of **0.857** on the held-out test set. Class-wise performance is strong for class 1 (precision **0.872**, recall **0.965**, $F_1 = \mathbf{0.916}$) and comparable to the best single models; for class 0, the ensemble yields precision 0.727, recall 0.400, $F_1 = 0.516$. Relative to the optimized MLP (accuracy 0.85, class 0 $F_1 = 0.53$), the ensemble offers a *slightly higher overall accuracy* and class 1 $F_1$ while maintaining a similar class 0 $F_1$, but lower recall for class 0. Given our clinical emphasis on identifying class 0, we retain the MLP as the primary model and use the ensemble as a robustness check.

### 7.0.10. AdaBoost ensemble

We also evaluated *AdaBoost*, which builds a stage-wise additive model of weak learners by reweighting training examples according to previous errors.

Starting from uniform weights, each round fits a base classifier (here, decision stumps/trees); misclassified samples receive larger weights and correctly classified samples smaller weights, producing a new learner that focuses on the hardest cases. The final prediction is a weighted vote of all stages, with weights proportional to each learner's accuracy on reweighted data.

As shown in Fig. 12, AdaBoost achieved overall accuracy **0.8476** on the held-out test set. Class-wise performance was strong for class 1 (precision 0.8791, recall 0.9412, $F_1 = 0.9091$) and moderate for class 0 (precision 0.6429, recall 0.4500, $F_1 = 0.5294$). Although AdaBoost is competitive overall, the optimized MLP remains our preferred model because it attains a comparable accuracy (0.85) while offering the strongest class 0 $F_1$ among the candidates, which aligns with our clinical objective of flagging likely non-responders.

### 7.0.11. Stacking ensemble

We constructed a *stacked generalization* (stacking) model that learns how to combine the out-of-fold predictions from multiple base learners into a final meta-prediction. The base layer comprised Logistic Regression, SVM, Random Forest, Naïve Bayes, MLP, and XGBoost, each tuned on the training split. For stacking, we generated out-of-fold predictions from each base model using stratified $k$-fold CV on the training data and used these as features for a meta-learner (logistic regression), which was then trained on the full out-of-fold matrix and applied to the held-out test set. This design reduces information leakage and lets the meta-learner exploit complementary error patterns across models.

As shown in Fig. 13, the stacking ensemble achieved accuracy **0.8095**. Class 1 performance remained strong (precision 0.8652, recall 0.9059, $F_1 = 0.8851$), while class 0 performance was moderate (precision 0.5000, recall 0.4000, $F_1 = 0.4444$). Compared with the optimized MLP (accuracy 0.85, class 0 $F_1 = 0.53$), stacking did not improve minority-class detection and thus is not preferred for our clinical objective, though it provides a useful robustness baseline.

### 7.0.12. Can foundation models help? (TabPFN)

We explored the *TabPFN* foundation model for tabular classification, which uses a Transformer pre-trained to approximate Bayesian posterior predictive inference on synthetic tasks and then performs *inference-only* fine-tuning at test time via forward passes (no gradient updates) [53]. On our

held-out test set, TabPFN achieved accuracy **0.83** (Fig. 14). Class 1 performance remained high (precision 0.83, recall 0.99, $F_1 = 0.90$), but class 0 recall was low (0.15; $F_1 = 0.25$). Given our clinical objective to *identify likely non-responders* (class 0), the optimized MLP (class 0 $F_1 = 0.53$) remains preferable despite TabPFN's strong overall discrimination for class 1.

*7.0.13. Benchmarking against human experts and a generative AI baseline*

To validate clinical utility and assess reliability, we benchmarked our model against *human experts* and a *generative AI* baseline. From the held-out test set we constructed an evaluation subset of $n = 30$ anonymized CRS cases. We first obtained class probabilities from the optimized MLP, $\hat{p}_i = \Pr(y_i{=}1 \mid \mathbf{x}_i)$, and quantified prediction uncertainty as $|\hat{p}_i - 0.5|$. We then stratified by uncertainty to obtain a balanced difficulty mix: the 10 *most uncertain* cases (smallest $|\hat{p}_i - 0.5|$), 10 *most certain* cases (largest $|\hat{p}_i - 0.5|$), and 10 *moderately certain* cases (middle of the remaining list). The final set preserved the approximate class ratio of the test split and was used for blinded comparison.

*Human and LLM baselines.* Six board-certified Otolaryngology physicians who practice as subspecialty Rhinologists independently labeled each case as *Surgery expected to result in desired improvement in QoL (1)* vs. *Surgery will not result in desired improvement in QoL (0)* using the same tabular summary provided to the model. Surgeons were selected based on participation in the original outcomes study to confirm domain knowledge on CRS surgical outcomes and to ensure familiarity with the clinical data fields. To facilitate expert labeling at scale, we deployed a secure web application for case review and adjudication (`http://149.166.246.230:4000/patientForms/`). The tool supports authenticated access, session persistence (save-and-resume), and full edit history so clinicians can revise prior responses. A custom graphical user interface (GUI) was developed to collect outcome predictions from these six experienced Rhinologists with expertise in CRS outcomes, all of whom were familiar with the study design and data format. Using the GUI, each expert reviewed detailed subject-level metadata and classified the expected treatment outcome as either "surgical success" or "surgical failure." Experts also rated their confidence on a five-point Likert scale ranging from "very confident" to "not at all confident." Cases were presented in a standardized tabular format with embedded guidance, and submissions were time-stamped and stored server-side to enable auditability and inter-rater

analyses. In parallel, we evaluated a large language model (ChatGPT, GPT-5 Thinking) using prompt-engineering to elicit a binary decision from the same structured inputs.[1]

*Results.* Figure 15 shows confusion matrices for our MLP and ChatGPT on the 30 cases. The MLP achieved (accuracy = 80%), with class-wise F1 metrics of $F1_0 = 0.57$, $F1_1 = 0.87$, while ChatGPT obtained accuracy = 57%, $F1_0 = 0.38$, $F1_1 = 0.67$. Notably, both methods recovered a similar class 0 recall (0.57), but ChatGPT exhibited substantially lower class 1 recall (0.57 vs. 0.87), producing many false negatives (under-calling surgery) in this cohort.

Table 3 summarizes accuracies for the expert predictions, our trained MLP and the ChatGPT model. The MLP outperforms the ChatGPT model and performs on par or better than the expert physicians on this subset (the average physician accuracy being 75.6%), highlighting the model's potential as a decision-support tool. Importantly, the uncertainty-stratified sampling increases face validity of the comparison by ensuring a mix of easy, ambiguous, and hard cases rather than only "cherry-picked" examples.

In addition to reporting individual expert performance, we also formed a *panel-of-experts* prediction for each case by aggregating the six independent expert labels via majority vote. In the event of an exact tie (3 vs. 3), we used the experts' self-reported confidence as a tie-breaker by selecting the class with the larger summed confidence across experts voting for that class. This aggregation yields a single consensus decision per case and approximates how multidisciplinary input might be combined in practice. Table 4 summarizes per-expert metrics and the pooled panel performance on the 30-case benchmarking subset described above. The full six-expert panel achieved 76.7% accuracy (23/30) with two ties resolved by confidence weighting. As a sensitivity analysis, we also evaluated a reduced panel consisting of the two highest-accuracy experts on this subset. We again resolved disagreements using the confidence-weighted tie-break rule described above. The resulting top-2 panel performance is reported in Table 4. Notably, panel aggregation did not improve accuracy relative to the best individual experts, likely because errors were correlated across raters with higher mispredictions, so majority voting gets dominated by shared misinterpretations rather than

---

[1]We did not fine-tune the LLM; prompts asked for a deterministic decision and brief rationale.

Table 3: Thirty-case human benchmarking subset: **accuracy** by rater. Row 1 reports Doctors 1–4; Row 2 reports Doctors 5–6, the in-house MLP, and ChatGPT.

|                | Doctor 1 | Doctor 2 | Doctor 3 | Doctor 4 |
|----------------|----------|----------|----------|----------|
| **Accuracy (%)** | 76.67    | 66.70    | 70.00    | 83.30    |

|                | Doctor 5 | Doctor 6 | MLP (ours) | ChatGPT |
|----------------|----------|----------|------------|---------|
| **Accuracy (%)** | 80.00    | 76.67    | 80.00      | 56.67   |

Table 4: Expert benchmarking on the 30-case subset: per-expert performance and panel aggregation. $P_0$ and $P_1$ represent precision for class 0 and 1, respectively, and similarly $R_0$ and $R_1$ represent recall for class 0 and 1, respectively. "Panel of experts" uses majority vote; ties are resolved using confidence-weighted voting.

| Rater | Acc | $P_0$ | $R_0$ | $P_1$ | $R_1$ |
|-------|-----|-------|-------|-------|-------|
| Doctor 1 | 0.767 | 0.500 | 0.286 | 0.808 | 0.913 |
| Doctor 2 | 0.667 | 0.286 | 0.286 | 0.783 | 0.783 |
| Doctor 3 | 0.700 | 0.333 | 0.286 | 0.792 | 0.826 |
| Doctor 4 | 0.833 | 0.750 | 0.429 | 0.846 | 0.957 |
| Doctor 5 | 0.800 | 0.667 | 0.286 | 0.815 | 0.957 |
| Doctor 6 | 0.767 | 0.500 | 0.143 | 0.786 | 0.957 |
| Panel of experts (6) | 0.767 | 0.500 | 0.143 | 0.786 | 0.957 |
| Top-2 panel (Doctors 4 and 5) | 0.800 | 0.667 | 0.286 | 0.815 | 0.957 |

averaging out independent noise.

*Performance gradient across case difficulty (Hard → Medium → Easy).* Because the 30-case benchmark set was constructed by stratifying cases according to the MLP's prediction uncertainty, it naturally supports an analysis of whether both human and model performance varies across difficulty tiers. When pooling expert judgments within each tier, accuracy increased monotonically from Hard to Easy (average accuracy - Hard: 0.55; Medium: 0.783; Easy: 0.933), indicating a clear difficulty gradient (Fig. 16). Figure 17 further shows that most individual experts exhibit the same Hard→Medium→Easy pattern, with the pooled mean tracking this monotonic rise.

The optimized MLP exhibited a similarly monotonic gradient across the same tiers, achieving accuracies of 0.60 (Hard), 0.80 (Medium), and 1.00 (Easy). Together, these results suggest that both clinician and model er-

rors concentrate in the most ambiguous (Hard) cases, while performance approaches ceiling on highly certain (Easy) cases.

In terms of comparison with the generative AI model, across the uncertainty spectrum, our MLP maintained balanced performance and substantially higher sensitivity for class 1, while ChatGPT tended to under-call surgery (false negatives). With comparable performance to physicians in terms of prediction accuracy, the MLP's +23.3 percentage point advantage over ChatGPT (80.0% vs. 56.7%) and its strong class-wise F1 scores further support its applicability as a complementary decision-support tool to aid clinicians, rather than a replacement for expert judgment.

### 7.0.14. Feature importance analysis

Understanding which variables drive model predictions is essential for clinical adoption and for generating testable hypotheses. We therefore quantified *global* feature importance for the optimized MLP using **permutation importance** on the held-out test set, scored by the drop in balanced accuracy ($\Delta$ BA) after randomly permuting each feature while holding all others fixed. To reduce Monte Carlo noise, we averaged importance over multiple permutations per feature and report the mean $\pm$ standard deviation as error bars (Fig. 18). This procedure probes the model's *reliance* on each input without assuming linearity and is invariant to feature scaling.

*Key findings.* The baseline SNOT-22 total score (`SNOT22_BLN_TOTAL`) was by far the most influential predictor, producing the largest decrease in balanced accuracy when permuted. This aligns with clinical intuition: preoperative symptom burden is strongly associated with postoperative patient-reported outcomes and thus with the model's recommendation signal. Age and imaging disease burden (baseline CT score, `BLN_CT_TOTAL`) were the next most important, consistent with known associations between disease severity, remodeling, and treatment response. Positive allergy testing and prior surgery also ranked highly, suggesting that atopy and surgical history modulate expected benefit. Additional contributors included presence of nasal polyps (`CRS_POLYPS`), COPD, and septal deviation (`SEPT_DEV`), all biologically plausible correlates of disease complexity and airflow/ventilation. Socioeconomic proxies (`HOUSEHOLD_INCOME`, `INSURANCE`) and demographics (`SEX`, `RACE`) showed modest but nonzero importance, indicating potential potential social determinants of health patterns; these variables should be monitored for causality rather than interpreted causally. Lower-ranked co-

morbidities (e.g., diabetes, OSA, asthma) contributed incrementally, consistent with their secondary roles in CRS symptomatology.

*Clinical plausibility and caveats.* Overall, the ranking mirrors domain expectations: symptoms $\rightarrow$ anatomic/imaging severity $\rightarrow$ atopy/surgical history $\rightarrow$ comorbidity/socioeconomic context. Because permutation importance reflects the model's *predictive dependence* and not causal effects, correlated features can share or mask importance, and interactions may distribute influence across variables. Accordingly, we treat these results as a global sanity check and pair them with local explanations (e.g., SHAP) next. This combination improves interpretability and supports responsible deployment in decision-support workflows.

### 7.0.15. Global and local interpretability via SHAP

To complement permutation importance and provide *directional*, patient-level explanations, we computed Shapley additive explanations (SHAP) for the optimized MLP. SHAP builds on the Shapley value from cooperative game theory [54] and provides consistent, locally faithful attributions for ML models [55]. We used the standard Kernel/Deep SHAP interface on the held-out test set and summarize results in Fig. 19.

*Global ranking and directionality.* The mean-|SHAP| bar plot (left) broadly corroborates the permutation ranking: the SNOT-22 baseline score dominates, followed by household income, baseline endoscopy score, positive allergy testing, baseline CT score, and age, with anatomic factors (e.g., septal deviation), history (e.g., previous surgery), and polyps contributing next. The beeswarm plot (right) adds *sign*: high preoperative symptom burden (shown with SNOT-22 baseline score, red points) concentrates on the positive SHAP side, pushing the model toward $\hat{y}=1$ (surgery), whereas low values (blue) push toward $\hat{y}=0$. A similar trend is visible for objective disease measures (baseline endoscopy score, baseline CT score) and atopy (positive allergy testing); higher values generally increase the predicted probability of surgery, consistent with domain expectations that greater symptom and anatomic burden favor operative management.

*Heterogeneity and fairness considerations.* Socioeconomic variables (like household income, insurance) exhibit substantial spread in SHAP values, indicating heterogeneous effects across patients. These social factors are important to surgical outcomes as they determine care access, medicine availability, diet,

and air quality. While these features improve discrimination, they can also encode access or practice patterns rather than biology; we therefore interpret them cautiously, monitor their contributions, and report fairness diagnostics separately. Demographics (sex, race) and several comorbidities (e.g., diabetes, smoker) show smaller average contributions, aligning with their lower global importances.

*Clinical plausibility and use.* Together, SHAP and permutation analyses produce a coherent narrative: symptoms and endoscopic/CT severity are primary drivers, modulated by atopy and surgical history, with comorbidities and socioeconomic context providing secondary signal. SHAP's local attributions are particularly useful for case review: clinicians can see *which* factors pushed an individual recommendation toward surgery vs. no surgery, supporting transparent, auditable decision support. These also aid in interpretability in a real-time clinical decision support application in community health settings where there is limited time and perhaps the doctor and the patient are not particularly familiar with data science techniques.

## 8. Discussion

This work demonstrates that routinely available preoperative data can support clinically useful individualized prediction of endoscopic sinus surgery (ESS) outcomes in Chronic Rhinosinusitis (CRS). Among multiple supervised learners, an optimized multilayer perceptron (MLP) provided the best balance of overall accuracy ($\approx 0.85$ on the held-out test set) and minority-class detection (class 0 F1 $\approx 0.53$), outperforming classic baselines (logistic regression, SVM, naive Bayes, random forest) and more sophisticated alternatives (XGBoost, stacking/boosting ensembles). Gains from ensembling were modest, consistent with a regime in which model variance is already low relative to the noise and sample size constraints of the problem. Human benchmarking on a stratified 30-case subset further showed that the MLP exceeded a strong generative-AI baseline and performed comparably to expert physicians, indicating that signal present in standard clinical variables is sufficient to aid decision support. Additionally, the human expert feedback provides valuable clinician intuition and enables comparing the factors deemed as salient by the ML model and the expert.

*Interpretability and clinical plausibility.* Interpretability was a desired property of our modeling. Global permutation importance and SHAP analyses

29

converged on a clinically coherent hierarchy: disease severity measures including patient reported SNOT-22, objective clinician scores for CT and endoscopy exams, and atopy/surgical history were the principal drivers of the model's recommendations, with comorbidities and socioeconomic context providing secondary signal. There can be a complex and individualized interplay among these many factors, the exploration of which could be an interesting future work.

*Class imbalance and error asymmetry.* Across all algorithms, performance on class 0 (patients unlikely to achieve the desired outcome) trailed class 1. This asymmetry is expected: class 0 is less frequent and potentially more heterogeneous. In many medical outcome settings, minority classes represent non-response or adverse events that are under-sampled, noisier, and driven by unmeasured factors. Our results reflect this: while class 1 F1 $\approx$ 0.91, class 0 F1 $\approx$ 0.53 for the MLP. We tried to mitigated the imbalance through class weighting, but a fundamental data limitation remains. Future work will explicitly address imbalance with re-sampling (oversampling of the minority class) and synthetic data generation (e.g., SMOTE/ADASYN, mixup); and model-based augmentation (e.g., VAEs/GANs) constrained to clinically plausible regions of feature space. Active learning with expert feedback focused on uncertain/rare cases could further enrich the minority class efficiently.

*Limits of the current study.* Several considerations exist, of course. First, the dataset is moderate in size. Although we used conservative preprocessing and assessed robustness, sample size likely caps achievable accuracy and minority recall. Second, our endpoint is binarized at six months. Some non-responders may benefit later, and continuous or longitudinal outcomes (e.g., $\Delta$SNOT-22 trajectories) could be more informative, although our prior study shows this is limited [56]. Third, socioeconomic variables improved prediction but can encode access or practice patterns. As such, they should be treated as potential confounders and subgroup performance should be monitored for fairness, which we plan to do in the future. Fourth, while the human benchmarking is encouraging, it was conducted on a small stratified subset of samples against legitimate experts in the field. We do not know how the average physicians in the community will compare, who have different experiences and may not keep up with the nuances in the literature or dive deeply into such conversations. In terms of the generative AI comparison, we experimented with a single LLM configuration. Overall, broader reader studies and varied

LLM baselines are warranted. Finally, external generalizability remains to be proven; site effects, referral patterns, surgical technique and perioperative care, medication adherence, patient psychology, priming effects for survey responses are all considerations that may shift feature distributions.

*Clinical relevance.* Even with these constraints, $\approx 85\%$ accuracy with transparent explanations is meaningful for preoperative counseling. Errors are not equivalent clinically: false negatives (predicting non-benefit when benefit is likely) may delay effective surgery, while false positives may expose a patient to operative risk with limited expected gain. Our framework supports threshold tuning and individualized decision curves, enabling clinicians to select operating points aligned with patient preferences and risk tolerance. The interactive visualizations help communicate why the recommendation was made and how it might change under alternative assumptions.

*Future directions.* We outline several priorities to translate these findings:

1. **Data growth and harmonization:** expand to broader cohorts; standardize variable definitions; leverage federated learning for privacy-preserving training. Also, include other surgery specific and patient psychological factors, other outcomes of relevance for decision making.
2. **Learning under imbalance:** integrate calibrated re-weighting, adopt techniques like SMOTE/ADASYN/mixup, and generative augmentation with clinical plausibility checks; evaluate minority-aware metrics prospectively.
3. **Temporal/external validation:** perform temporal splits, external site validation, and pre-registered prospective studies; assess domain shift and transportability.
4. **Better targets and calibration:** model continuous and time-to-event endpoints; report Brier score/ECE; apply temperature/Platt calibration; use decision-curve analysis to quantify net benefit across thresholds.
5. **Uncertainty and reliability:** add Bayesian/post-hoc uncertainty quantification, conformal prediction, and reject-option policies to flag cases for multidisciplinary review.
6. **Individualized treatment benefit:** move beyond outcome prediction to estimate conditional treatment effects (uplift models, causal forests) to quantify patient-specific expected gain from surgery versus medical therapy.

7. **Human–AI collaboration:** expand reader studies, measure assistance effects (with/without model), and optimize the UI for rapid, explainable triage; incorporate clinician feedback loops for continual learning.

8. **Foundation/self-supervised models:** our TabPFN exploration showed strong class 1 discrimination but limited minority recall; future work will test tabular foundation models and hybrid pipelines (self-supervised feature learning + calibrated discriminators) under imbalance constraints.

## 9. Conclusion

We present an interpretable, data-efficient pipeline for predicting CRS surgical outcomes. An optimized MLP achieved strong discrimination on a held-out cohort and surpassed both a generative-AI baseline and expected human performance on a stratified reader set, while providing transparent global and local explanations. The principal limitation is minority-class performance driven by class imbalance and sample size, ubiquitous challenges in medical outcome prediction. Addressing these with targeted data growth, imbalance-aware learning, rigorous external validation, and careful calibration will be key to translating such models from retrospective studies to prospective, decision-support tools that improve patient counseling and, ultimately, outcomes.

## ACKNOWLEDGMENT

## Appendix A. Tripod AI Check

I checked mostly, we may need to confirm 1 point: fairness.

# References

[1] P. Malik, M. Pathania, and V. K. Rathaur, "Overview of artificial intelligence in medicine," *Journal of Family Medicine and Primary Care*, vol. 8, no. 7, pp. 2328–2331, 2019.

[2] Y. Mintz and R. Brodie, "Introduction to artificial intelligence in medicine," *Minimally Invasive Therapy & Allied Technologies*, vol. 28, no. 2, pp. 73–81, 2019.

[3] V. H. Buch, I. Ahmed, and M. Maruthappu, "Artificial intelligence in medicine: Current trends and future possibilities," *British Journal of General Practice*, vol. 68, no. 668, pp. 143–144, 2018.

[4] N. R. Sahni and B. Carrus, "Artificial intelligence in us health care delivery," *New England Journal of Medicine*, vol. 389, no. 4, pp. 348–358, 2023.

[5] C. Varghese, E. M. Harrison, G. O'Grady, and E. J. Topol, "Artificial intelligence in surgery," *Nature medicine*, vol. 30, no. 5, pp. 1257–1268, 2024.

[6] K.-H. Yu, A. L. Beam, and I. S. Kohane, "Artificial intelligence in healthcare," *Nature biomedical engineering*, vol. 2, no. 10, pp. 719–731, 2018.

[7] H.-J. Suh, J. Son, and K. Kang, "Application of artificial intelligence in the practice of medicine," *Applied Sciences*, vol. 12, no. 9, p. 4649, 2022.

[8] F. S. Collins and H. Varmus, "A new initiative on precision medicine," *New England journal of medicine*, vol. 372, no. 9, pp. 793–795, 2015.

[9] P. W. Hellings, W. J. Fokkens, C. Bachert, C. A. Akdis, T. Bieber, I. Agache, M. Bernal-Sprekelsen, G. W. Canonica, P. Gevaert, G. Joos, *et al.*, "Positioning the principles of precision medicine in care pathways for allergic rhinitis and chronic rhinosinusitis–a euforea-aria-eposairways icp statement," *Allergy*, vol. 72, no. 9, pp. 1297–1305, 2017.

[10] A. Muraro, R. F. Lemanske Jr, P. W. Hellings, C. A. Akdis, T. Bieber, T. B. Casale, M. Jutel, P. Y. Ong, L. K. Poulsen, P. Schmid-Grendelmeier, *et al.*, "Precision medicine in patients with allergic diseases: airway diseases and atopic dermatitis—practall document of the

european academy of allergy and clinical immunology and the american academy of allergy, asthma & immunology," *Journal of allergy and clinical immunology*, vol. 137, no. 5, pp. 1347–1358, 2016.

[11] T. J. Loftus, P. J. Tighe, A. C. Filiberto, *et al.*, "Artificial intelligence and surgical decision-making," *JAMA Surgery*, vol. 155, no. 2, pp. 148–158, 2020.

[12] L. Rudmik, "Economics of chronic rhinosinusitis," *Current Allergy and Asthma Reports*, vol. 17, no. 4, p. 20, 2017.

[13] R. R. Orlandi, T. T. Kingdom, T. L. Smith, *et al.*, "International consensus statement on allergy and rhinology: Rhinosinusitis 2021," *International Forum of Allergy & Rhinology*, vol. 11, no. 3, pp. 213–739, 2021.

[14] P. T. Le, Z. M. Soler, R. Jones, J. L. Mattos, S. A. Nguyen, and R. J. Schlosser, "Systematic review and meta-analysis of snot-22 outcomes after surgery for chronic rhinosinusitis with nasal polyposis," *Otolaryngology–Head and Neck Surgery*, vol. 159, no. 3, pp. 414–423, 2018.

[15] Z. M. Soler, R. Jones, P. Le, L. Rudmik, J. L. Mattos, S. A. Nguyen, and R. J. Schlosser, "Sino-nasal outcome test-22 outcomes after sinus surgery: A systematic review and meta-analysis," *The Laryngoscope*, vol. 128, no. 3, pp. 581–592, 2018.

[16] C. Hopkins, R. Hettige, A. Soni-Jaiswal, *et al.*, "Chronic rhinosinusitis outcome measures (chrome), developing a core outcome set for trials of interventions in chronic rhinosinusitis," *Rhinology*, vol. 56, no. 1, pp. 22–32, 2018.

[17] L. Rudmik and Z. M. Soler, "Medical therapies for adult chronic sinusitis: a systematic review," *Jama*, vol. 314, no. 9, pp. 926–939, 2015.

[18] W. W. Stevens, R. J. Lee, R. P. Schleimer, and N. A. Cohen, "Chronic rhinosinusitis pathogenesis," *Journal of Allergy and Clinical Immunology*, vol. 136, no. 6, pp. 1442–1453, 2015.

[19] A. Choi, S. Xu, A. U. Luong, and S. K. Wise, "Current review of comorbidities in chronic rhinosinusitis," *Current Allergy and Asthma Reports*, vol. 25, no. 1, p. 4, 2024.

[20] L. Rudmik, T. L. Smith, R. J. Schlosser, *et al.*, "Productivity costs in patients with refractory chronic rhinosinusitis," *The Laryngoscope*, vol. 124, no. 9, pp. 2007–2012, 2014.

[21] V. C. Pandrangi, J. C. Mace, J.-H. Kim, *et al.*, "Work productivity and activity impairment in patients with chronic rhinosinusitis undergoing endoscopic sinus surgery—a prospective, multi-institutional study," *International Forum of Allergy & Rhinology*, vol. 13, no. 3, pp. 216–229, 2023.

[22] S. S. Smith, R. Kim, and R. Douglas, "Is there a role for antibiotics in the treatment of chronic rhinosinusitis?," *Journal of Allergy and Clinical Immunology*, vol. 149, no. 5, pp. 1504–1512, 2022.

[23] G. E. Davis, R. S. Zeiger, B. Emmanuel, *et al.*, "Systemic corticosteroid-related adverse outcomes and health care resource utilization and costs among patients with chronic rhinosinusitis with nasal polyposis," *Clinical Therapeutics*, vol. 44, no. 9, pp. 1187–1202, 2022.

[24] Z. M. Soler, E. Wittenberg, R. J. Schlosser, J. C. Mace, and T. L. Smith, "Health state utility values in patients undergoing endoscopic sinus surgery," *The Laryngoscope*, vol. 121, no. 12, pp. 2672–2678, 2011.

[25] A. W. Chow, M. S. Benninger, I. Brook, J. L. Brozek, E. J. C. Goldstein, L. A. Hicks, G. A. Pankey, M. Seleznick, G. Volturo, E. R. Wald, and T. M. File, "Idsa clinical practice guideline for acute bacterial rhinosinusitis in children and adults," *Clinical Infectious Diseases*, vol. 54, no. 8, pp. e72–e112, 2012.

[26] S. S. Smith, R. C. Kern, R. K. Chandra, B. K. Tan, and C. T. Evans, "National burden of antibiotic use for adult rhinosinusitis," *Journal of Allergy and Clinical Immunology*, vol. 132, no. 5, pp. 1230–1232, 2013.

[27] K. A. Smith, R. R. Orlandi, and L. Rudmik, "Cost of adult chronic rhinosinusitis: A systematic review," *The Laryngoscope*, vol. 125, no. 7, pp. 1547–1556, 2015.

[28] L. Rudmik, C. Hopkins, A. Peters, T. L. Smith, R. J. Schlosser, and Z. M. Soler, "Patient-reported outcome measures for adult chronic rhinosinusitis: a systematic review and quality assessment," *Journal of Allergy and Clinical Immunology*, vol. 136, no. 6, pp. 1532–1540, 2015.

[29] J. L. Mattos, L. Rudmik, R. J. Schlosser, T. L. Smith, J. C. Mace, J. Alt, and Z. M. Soler, "Symptom importance, patient expectations, and satisfaction in chronic rhinosinusitis," *International Forum of Allergy & Rhinology*, vol. 9, no. 6, pp. 593–600, 2019.

[30] J. J. Shin, AAO–HNSF Guideline Development Group, *et al.*, "Clinical practice guideline: Surgical management of chronic rhinosinusitis," *Otolaryngology–Head and Neck Surgery*, 2025. AAO–HNSF Clinical Practice Guideline.

[31] X.-r. Kang, B. Chen, Y.-s. Chen, *et al.*, "A prediction modeling based on snot-22 score for endoscopic nasal septoplasty: A retrospective study," *PeerJ*, vol. 8, p. e9890, 2020.

[32] A. Raghavan, E. Sage, M. Al-Ghezi, M. Aboueisha, I. Prohnitchi, J. P. Giliberto, I. Humphreys, A. Jafari, and W. M. Abuzeid, "Predicting surgical outcomes in chronic rhinosinusitis from preoperative patient data: A machine learning approach," in *International Forum of Allergy & Rhinology*, Wiley Online Library, 2025.

[33] J. F. Cohen and P. M. Bossuyt, "Tripod+ ai: an updated reporting guideline for clinical prediction models," 2024.

[34] Coalition for Health AI (CHAI), "Responsible ai guide (raig) and raig executive summary." https://www.chai.org/workgroup/responsible-ai/responsibleai-guide-raig-and-raig-executive-summary/. Accessed November 2025.

[35] G. S. Collins, J. B. Reitsma, D. G. Altman, and K. G. M. Moons, "Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (tripod): The tripod statement," *Annals of Internal Medicine*, vol. 162, no. 1, pp. 55–63, 2015.

[36] A. J. Vickers and E. B. Elkin, "Decision curve analysis: A novel method for evaluating prediction models," *Medical Decision Making*, vol. 26, no. 6, pp. 565–574, 2006.

[37] E. W. Steyerberg, *Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating.* Cham: Springer, 2 ed., 2019.

[38] M. Nagendran, Y. Chen, C. A. Lovejoy, *et al.*, "Artificial intelligence versus clinicians: systematic review of design, reporting standards, and claims of deep learning studies," *BMJ*, vol. 368, p. m689, 2020.

[39] J. Brazier, J. Roberts, and M. Deverill, "The estimation of a utility-based algorithm from the sf-36 health survey," *Journal of Health Economics*, vol. 21, no. 2, pp. 271–292, 2002.

[40] C. M. Bishop, *Pattern Recognition and Machine Learning.* New York: Springer, 2006.

[41] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction.* New York: Springer, 2 ed., 2009.

[42] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning.* Cambridge, MA: MIT Press, 2016.

[43] D. M. W. Powers, "Evaluation: From precision, recall and f-measure to roc, informedness, markedness & correlation," *Journal of Machine Learning Technologies*, vol. 2, no. 1, pp. 37–63, 2011.

[44] F. Pedregosa, G. Varoquaux, A. Gramfort, and et al., "Scikit-learn: Machine learning in python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

[45] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.

[46] C. M. Bishop, *Pattern Recognition and Machine Learning.* Springer, 2006.

[47] K. P. Murphy, *Machine Learning: A Probabilistic Perspective.* MIT Press, 2012.

[48] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.

[49] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, pp. 273–297, 1995.

[50] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.

[51] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 2001.

[52] D. H. Wolpert, "Stacked generalization," *Neural Networks*, vol. 5, no. 2, pp. 241–259, 1992.

[53] S. Müller, M. Feindt, M. Hörmann, N. Hollmann, *et al.*, "Tabpfn: A Transformer that solves small tabular classification problems in a second," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

[54] L. S. Shapley, "A value for n-person games," in *Contributions to the Theory of Games, Vol. II* (H. W. Kuhn and A. W. Tucker, eds.), vol. 28 of *Annals of Mathematics Studies*, pp. 307–317, Princeton University Press, 1953.

[55] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.

[56] A. S. DeConde, J. C. Mace, J. A. Alt, L. Rudmik, Z. M. Soler, and T. L. Smith, "Longitudinal improvement and stability of the snot-22 survey in the evaluation of surgical management for chronic rhinosinusitis," in *International forum of allergy & rhinology*, vol. 5, pp. 233–239, Wiley Online Library, 2015.
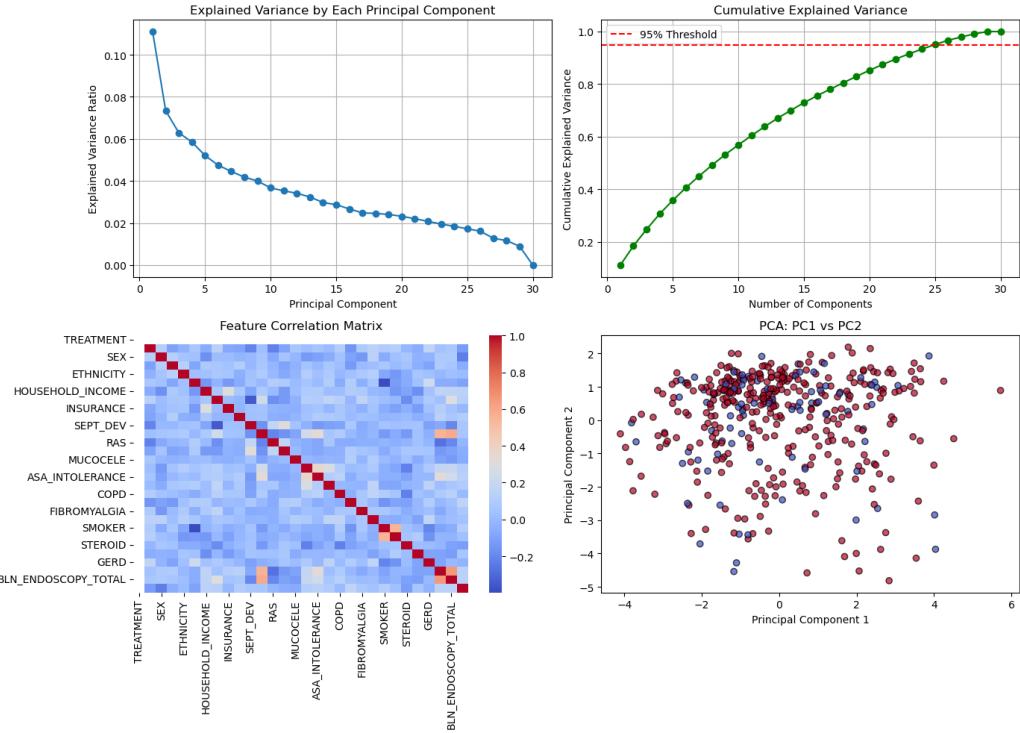
Figure 5: Exploratory feature and PCA analysis. (a) Explained variance by principal component. Each bar shows the proportion of total variance attributable to an individual component after standardization; early components carry the largest share, followed by a long tail. (b) Cumulative explained variance. The running total of variance captured as components are added; the dashed line marks 95% retained variance, indicating that a moderate subset of components suffices for compact representations. (c) Feature correlation heatmap (Pearson's correlation coefficient, $r$). Blue denotes negative and red positive associations among preprocessed clinical predictors; the matrix helps flag redundancy/collinearity before modeling. (d) PCA score plot (PC1 vs. PC2). Patients are projected onto the first two principal components and colored by class label; the overlap across classes in this linear 2D view suggests that discrimination likely depends on higher-order components and/or non-linear structure.

| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.75 | 0.30 | 0.43 | 20 |
| 1 | 0.86 | 0.98 | 0.91 | 85 |
| | | | | |
| Accuracy | | | 0.85 | 105 |
| | | | | |
| Macro Avg | 0.80 | 0.64 | 0.67 | 105 |
| Weighted Avg | 0.84 | 0.85 | 0.82 | 105 |



SVM on Held-Out Test Set: Classification Report and Confusion Matrix

| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.33 | 0.10 | 0.15 | 20 |
| 1 | 0.82 | 0.95 | 0.88 | 85 |
| | | | | |
| Accuracy | | | 0.79 | 105 |
| | | | | |
| Macro Avg | 0.58 | 0.53 | 0.52 | 105 |
| Weighted Avg | 0.73 | 0.79 | 0.74 | 105 |



Naïve Bayes on Held-Out Test Set: Classification Report and Confusion Matrix

| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.21 | 1.00 | 0.35 | 20 |
| 1 | 1.00 | 0.13 | 0.23 | 85 |
| | | | | |
| Accuracy | | | 0.30 | 105 |
| | | | | |
| Macro Avg | 0.61 | 0.56 | 0.29 | 105 |
| Weighted Avg | 0.85 | 0.30 | 0.25 | 105 |



MLP on Held-Out Test Set: Classification Report and Confusion Matrix

| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.64 | 0.45 | 0.53 | 20 |
| 1 | 0.88 | 0.94 | 0.91 | 85 |
| | | | | |
| Accuracy | | | 0.85 | 105 |
| | | | | |
| Macro Avg | 0.76 | 0.70 | 0.72 | 105 |
| Weighted Avg | 0.83 | 0.85 | 0.84 | 105 |



Random Forest on Held-Out Test Set: Classification Report and Confusion Matrix

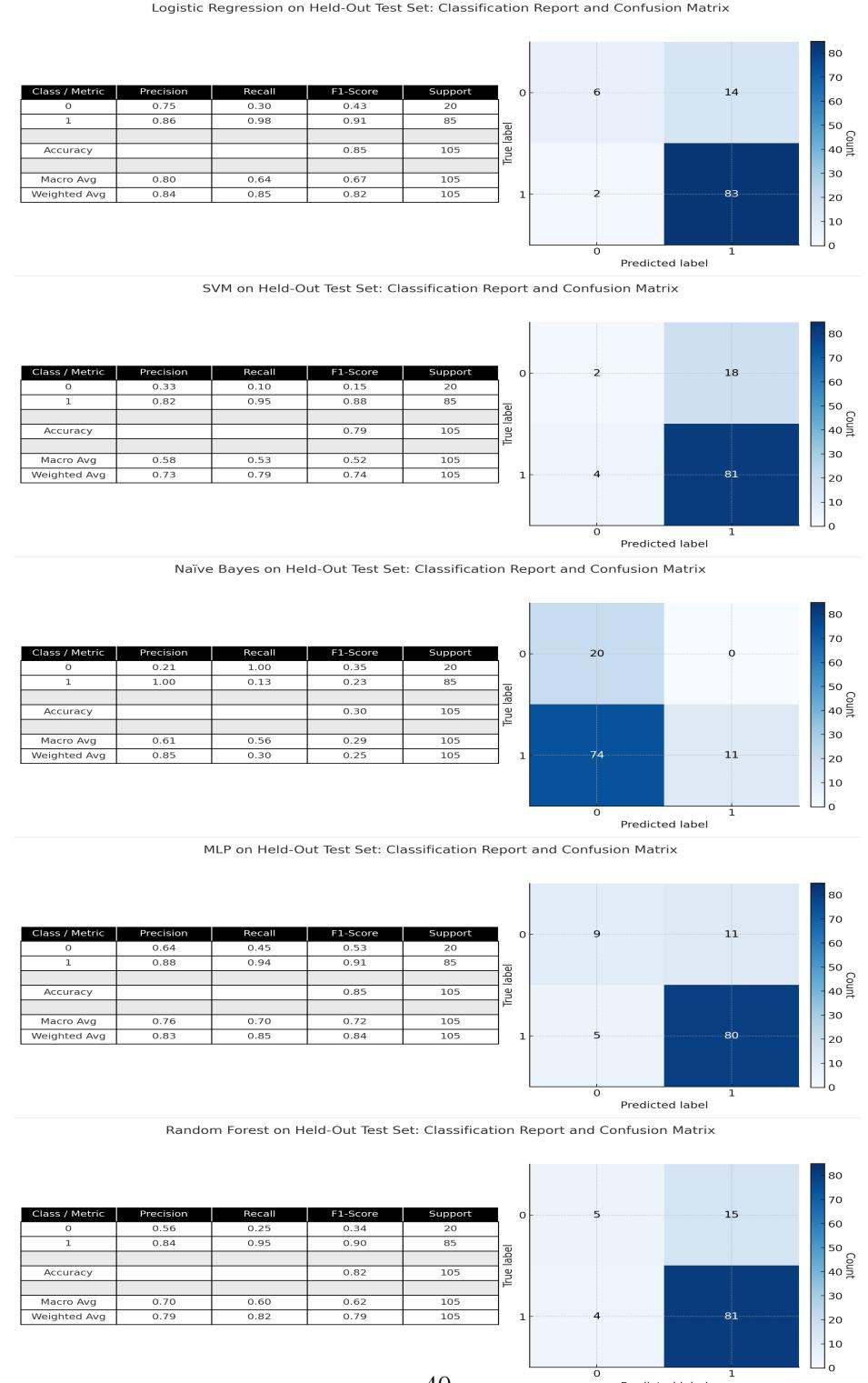| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.56 | 0.25 | 0.34 | 20 |
| 1 | 0.84 | 0.95 | 0.90 | 85 |
| | | | | |
| Accuracy | | | 0.82 | 105 |
| | | | | |
| Macro Avg | 0.70 | 0.60 | 0.62 | 105 |
| Weighted Avg | 0.79 | 0.82 | 0.79 | 105 |



40

Figure 6: Held-out test performance across five classifiers (Logistic Regression, SVM, Naïve Bayes, MLP, Random Forest). Each panel shows the classification report (precision, recall, F1-score, support) alongside the confusion matrix (counts).
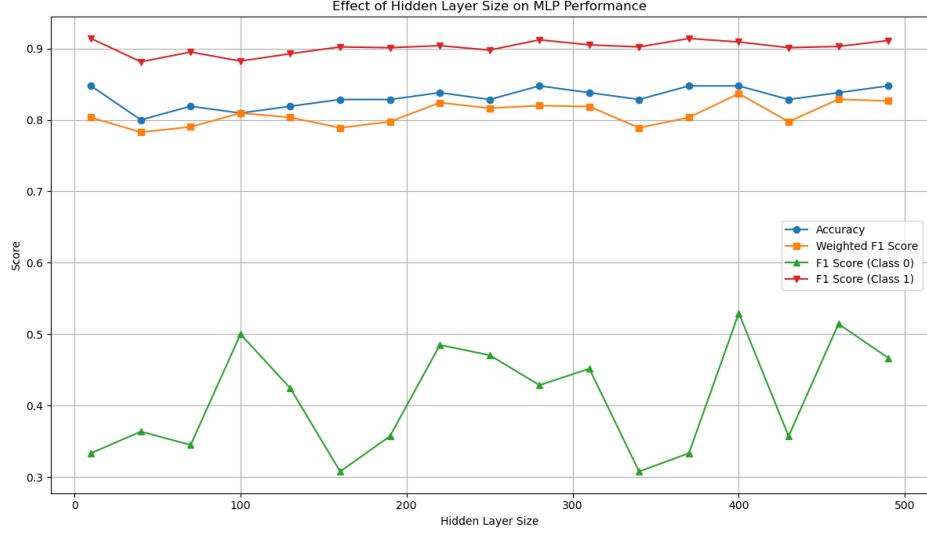
Figure 7: MLP hyperparameter tuning (single hidden layer). Validation performance vs. hidden-layer width showing accuracy, weighted F1, and class-specific F1 scores. Overall accuracy and class 1 F1 remain relatively stable across widths, whereas the class 0 F1 (our priority) improves with larger layers and peaks for wide settings ($\sim$400–480 units). Results are from stratified cross-validation on the training split.
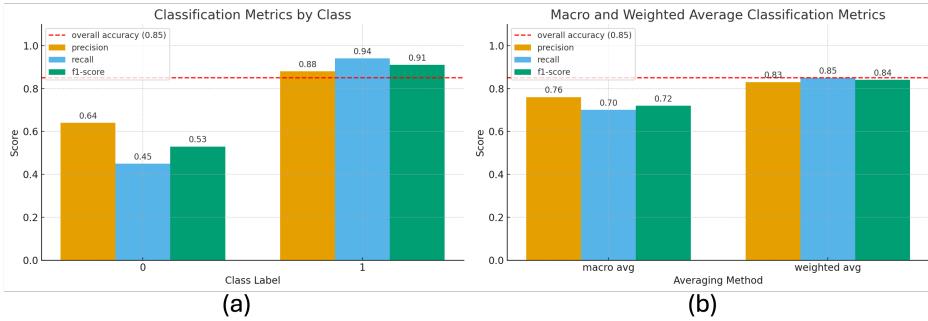


Figure 8: Optimized MLP performance on the held-out test set. Left: class-wise precision, recall, and F1 for class 0 and class 1 (dashed line marks overall accuracy, 0.85). Right: macro- and weighted-average metrics, which are close to overall accuracy, indicating stable performance without overfitting to the majority class.

Figure 9: Class-wise comparison of precision, recall, and F1-score across five models on the held-out test set (MLP, Random Forest, Logistic Regression, SVM, Naïve Bayes). Each panel shows grouped bars for class 0 vs. class 1. Overall, class 1 metrics are higher for all models, while the MLP attains the strongest class 0 F1 among the candidates.
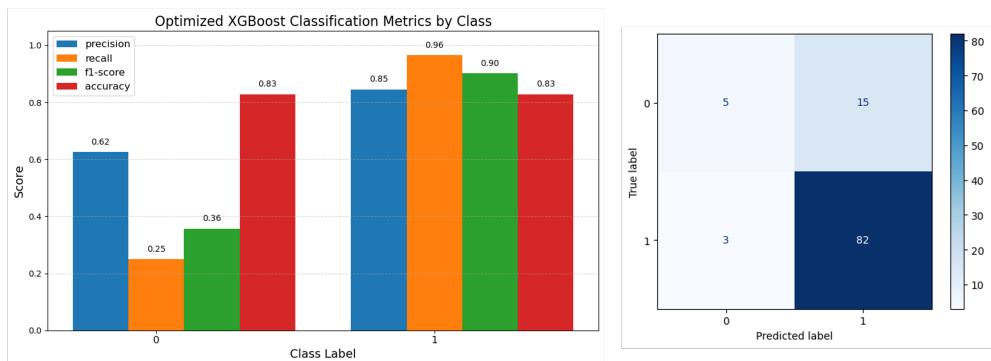


Figure 10: Optimized XGBoost results on the held-out test set. Left: class-wise precision, recall, $F_1$, and overall accuracy bars (class 0 vs. class 1). Right: confusion matrix (counts), yielding accuracy $\approx 0.83$ with strong class 1 recall (0.96) but modest class 0 recall (0.25).

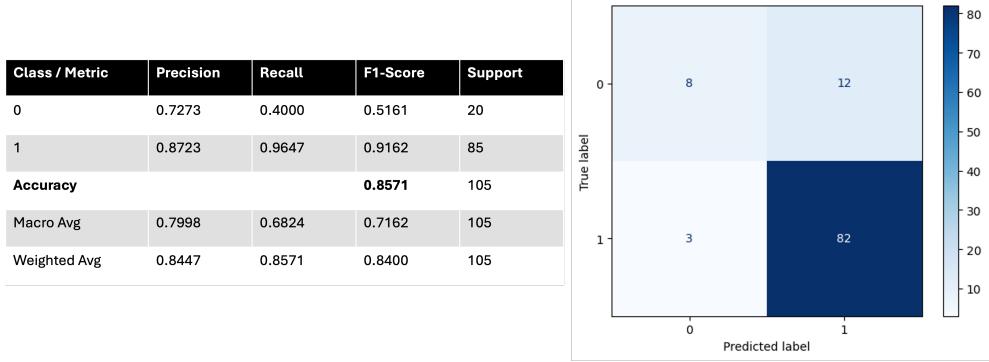| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.7273 | 0.4000 | 0.5161 | 20 |
| 1 | 0.8723 | 0.9647 | 0.9162 | 85 |
| **Accuracy** | | | **0.8571** | 105 |
| Macro Avg | 0.7998 | 0.6824 | 0.7162 | 105 |
| Weighted Avg | 0.8447 | 0.8571 | 0.8400 | 105 |

Figure 11: Majority-voting ensemble on the held-out test set. Left: classification report (precision, recall, $F_1$, support) for each class and macro/weighted averages. Right: confusion matrix (counts). The ensemble aggregates predictions from six tuned models (LR, SVM, RF, Naïve Bayes, MLP, XGBoost) with an MLP tie–break, yielding slightly higher overall accuracy than the single MLP while preserving similar class 0 performance.

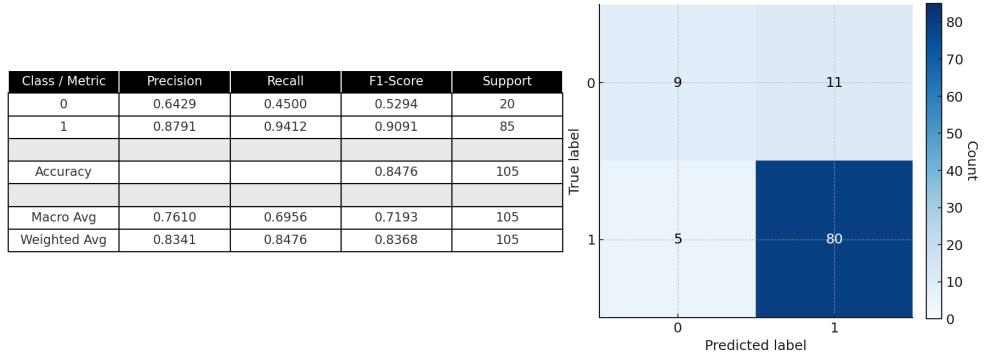AdaBoost Ensemble on Held-Out Test Set: Classification Report and Confusion Matrix

| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.6429 | 0.4500 | 0.5294 | 20 |
| 1 | 0.8791 | 0.9412 | 0.9091 | 85 |
| | | | | |
| Accuracy | | | 0.8476 | 105 |
| | | | | |
| Macro Avg | 0.7610 | 0.6956 | 0.7193 | 105 |
| Weighted Avg | 0.8341 | 0.8476 | 0.8368 | 105 |

Figure 12: AdaBoost ensemble on the held-out test set. Left: classification report (precision, recall, $F_1$, support) for each class and macro/weighted averages. Right: confusion matrix (counts) showing balanced overall performance with stronger class 1 detection and moderate class 0 recall.

43

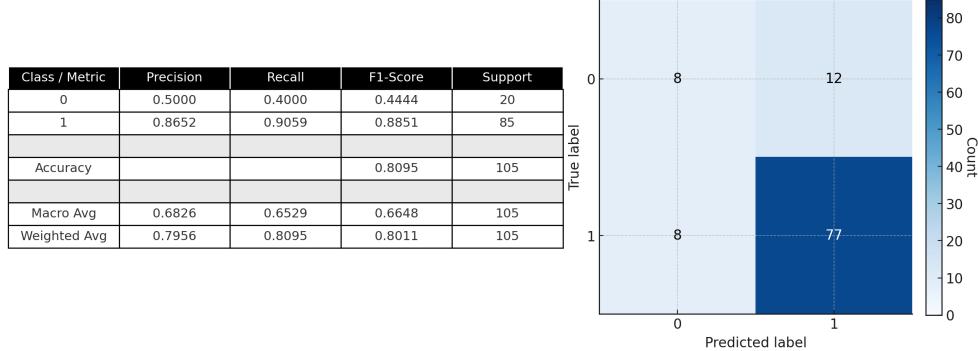Stacking Ensemble on Held-Out Test Set: Classification Report and Confusion Matrix

| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.5000 | 0.4000 | 0.4444 | 20 |
| 1 | 0.8652 | 0.9059 | 0.8851 | 85 |
| | | | | |
| Accuracy | | | 0.8095 | 105 |
| | | | | |
| Macro Avg | 0.6826 | 0.6529 | 0.6648 | 105 |
| Weighted Avg | 0.7956 | 0.8095 | 0.8011 | 105 |

Figure 13: Stacking ensemble on the held-out test set. Left: classification report by class and macro/weighted averages. Right: confusion matrix (counts). The meta-learner combines out-of-fold predictions from six tuned base models (LR, SVM, RF, Naïve Bayes, MLP, XGBoost); while overall accuracy is competitive, class 0 detection lags the MLP.

TabPFN on Held-Out Test Set: Classification Report and Confusion Matrix

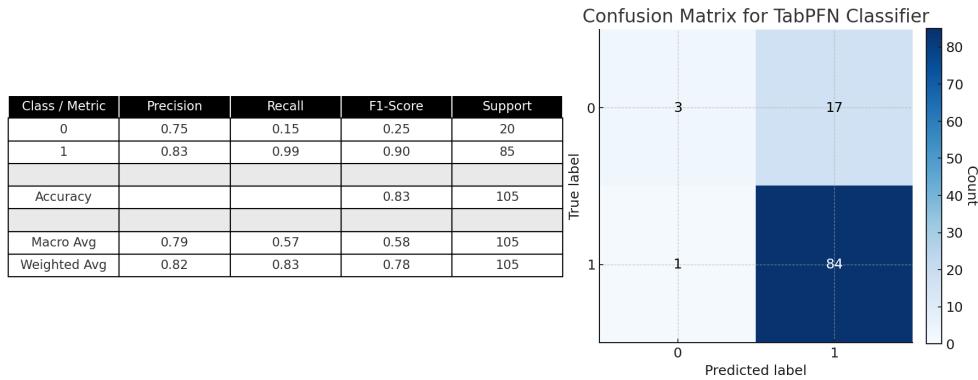| Class / Metric | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.75 | 0.15 | 0.25 | 20 |
| 1 | 0.83 | 0.99 | 0.90 | 85 |
| | | | | |
| Accuracy | | | 0.83 | 105 |
| | | | | |
| Macro Avg | 0.79 | 0.57 | 0.58 | 105 |
| Weighted Avg | 0.82 | 0.83 | 0.78 | 105 |

Figure 14: TabPFN on the held-out test set. Left: classification report by class and macro/weighted averages. Right: confusion matrix (counts). While TabPFN delivers excellent class 1 recall (0.99), class 0 recall (0.15) lags our optimized MLP, which is prioritized for clinical screening of likely non-responders.
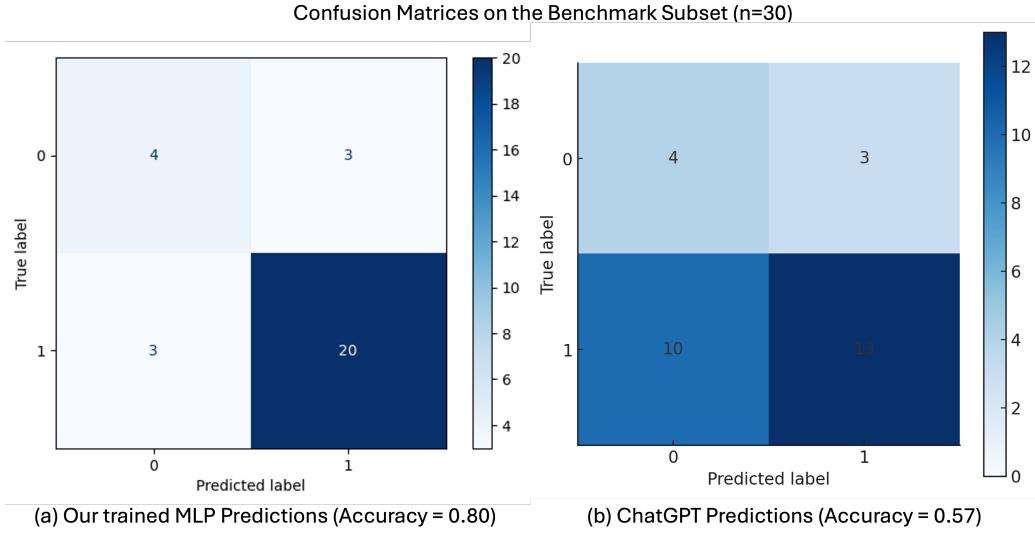
44

Confusion Matrices on the Benchmark Subset (n=30)

(a) Our trained MLP Predictions (Accuracy = 0.80)     (b) ChatGPT Predictions (Accuracy = 0.57)

Figure 15: Human-benchmark subset ($n = 30$). Left (a): MLP prediction confusion matrix (acc = 0.80. Right (b): ChatGPT prediction confusion matrix (acc = 0.57. The MLP shows markedly higher recall for class 1 while maintaining comparable recall for class 0, yielding superior overall discrimination on this stratified sample.
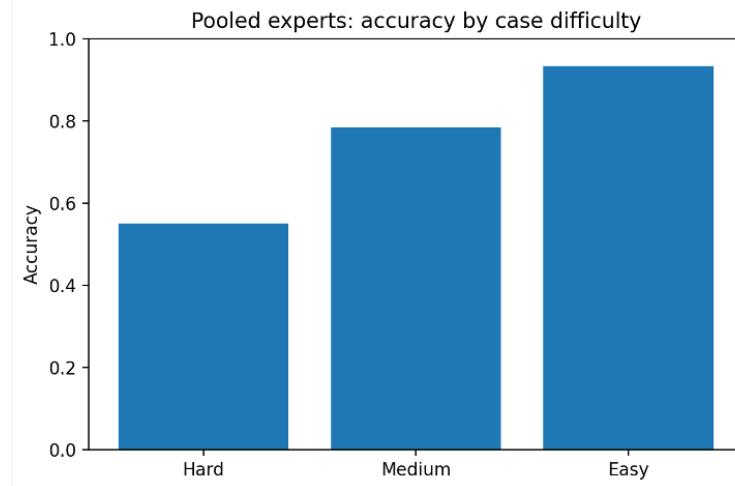


Figure 16: Pooled experts: accuracy by case difficulty tier (Hard/Medium/Easy) on the 30-case benchmark subset.
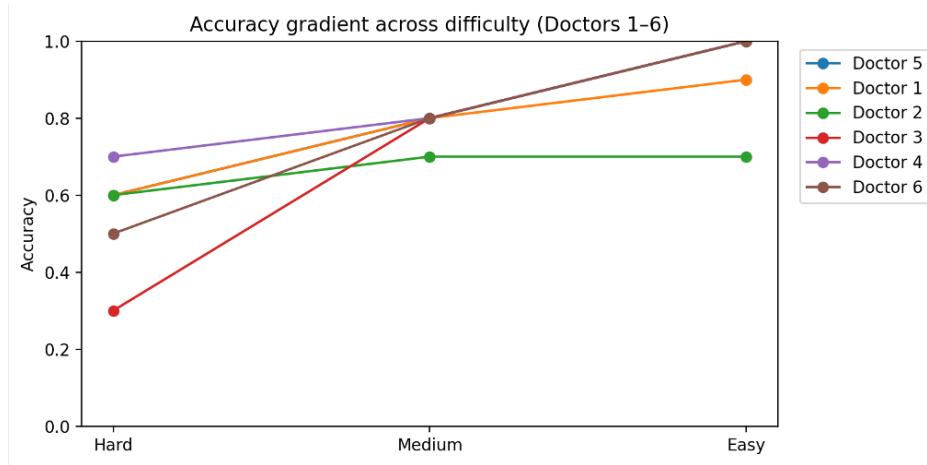
Figure 17: Accuracy gradient across difficulty tiers for each expert and the pooled mean. Most raters show a monotonic Hard→Medium→Easy increase.
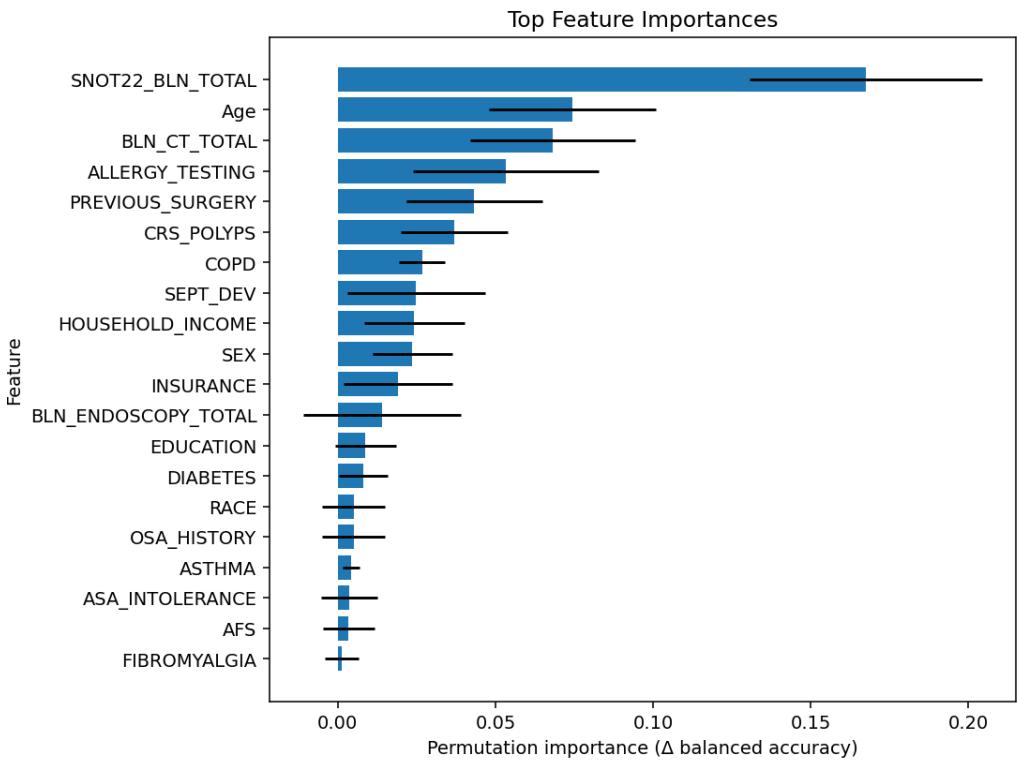
Figure 18: Permutation feature importance for the MLP classifier on the held-out test set. Bars show the mean decrease in balanced accuracy when each feature is randomly permuted (larger values indicate greater importance); black ticks denote the standard deviation across permutation repeats. The model is most sensitive to SNOT22_BLN_TOTAL, Age, and BLN_CT_TOTAL, followed by ALLERGY_TESTING, PREVIOUS_SURGERY, and CRS_POLYPS.
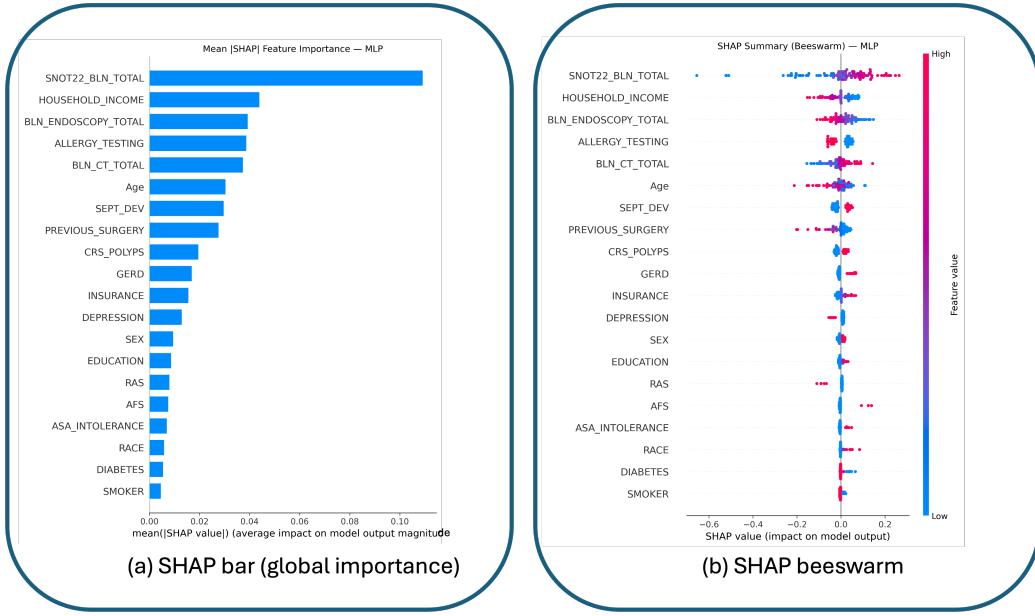
Figure 19: (a) SHAP bar (global importance): Mean absolute SHAP values on the test set rank features by their average contribution to the predicted probability of achieving the desired outcome; higher bars indicate greater global influence. (b) SHAP beeswarm: Each point represents a patient; the horizontal axis shows the SHAP value (feature effect on the log-odds/probability), and color encodes the feature value from low→high. Distributions reveal both effect size and heterogeneity.