

Predicting infected pancreatic necrosis in acute pancreatitis using machine learning models and feature selection

Received: 23 May 2025

Accepted: 29 January 2026

Published online: 28 February 2026

Cite this article as: Xin L., Yixuan D., Bohan H. *et al.* Predicting infected pancreatic necrosis in acute pancreatitis using machine learning models and feature selection. *Sci Rep* (2026). <https://doi.org/10.1038/s41598-026-38410-0>

Li Xin, Ding Yixuan, Huang Bohan, Shen Yunheng, Lv Hairong, Cao Feng, Yu Tong, Li Fei, Fei Xiaolu & Li Jia

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

Predicting Infected Pancreatic Necrosis in Acute Pancreatitis Using Machine Learning Models and Feature Selection

Li Xin^{1,†}, Ding Yixuan^{3,†}, Huang Bohan¹, Shen Yunheng², Lv Hairong²,
Cao Feng³, Yu Tong⁴, Li Fei^{3,*}, Fei Xiaolu^{3,*}, Li Jia^{3,4,*}

[†]Li Xin and Ding Yixuan contributed equally to this work.

*Correspondence:lifei@xwhospital.com; feixiaolu@xwhospital.com;

lijia@xwh.ccmu.edu.cn

Full list of author information is available at the end of the article.

Abstract

Background: Infected pancreatic necrosis (IPN) is a life-threatening complication of acute pancreatitis (AP), and its early prediction remains challenging. This study aimed to develop and externally validate interpretable machine learning models for individualized IPN risk prediction.

Methods: A total of 728 patients with AP admitted to Xuanwu Hospital, Capital Medical University, between 2017 and 2023 were retrospectively analyzed. Embedded feature selection was incorporated within model training using regularized linear and tree-based algorithms to enhance interpretability and prevent overfitting. Five machine learning algorithms and one neural network model were evaluated through nested cross-validation and an independent temporal external cohort consisting of 166 AP patients admitted to Xuanwu Hospital, Capital Medical University, between 2022 and 2023. Model discrimination, precision–recall, and probability calibration were assessed, and model explainability was analyzed using Shapley Additive Explanations (SHAP).

23 **Results:** The Random Forest model achieved the best overall performance, achieving
24 an external AUC of 0.764 (95% CI: 0.696–0.830, $P < 0.001$), precision of 0.893, recall of
25 0.604, and the lowest Brier score, indicating reliable probability calibration. SHAP analysis
26 identified Fibrinogen, APACHE II score, D-dimer, IL-6, and C-reactive protein (CRP) as
27 key predictors associated with increased IPN risk, while higher Lymphocyte count, and
28 Hematocrit were protective. These findings are consistent clinical pathophysiology.

29 **Conclusion:** The interpretable Random Forest model demonstrated robust discrimina-
30 tion and calibration for IPN prediction, providing a transparent and data-driven framework
31 for early risk stratification in acute pancreatitis. Prospective multicenter validation is war-
32 ranted before clinical implementation.

33 **Keywords**

34 Acute pancreatitis ; Deep learning ; Infected pancreatic necrosis ; Machine learning ; Prognosis
35 prediction

36 **Introduction**

37 Acute pancreatitis (AP) is a prevalent acute digestive system disorder characterized by its high
38 incidence and potential for severe outcomes. Globally, the incidence of AP ranges from 13
39 to 45 cases per 100,000 individuals annually, with a mortality rate of approximately 1%-3%
40 in mild cases and up to 20%-30% in severe cases [1, 2]. According to the modified Atlanta
41 classification, AP is categorized into mild, moderate, and severe based on the presence of organ
42 dysfunction and local complications, such as acute pancreatic pseudocysts, fluid collections,
43 and infected pancreatic necrosis (IPN) [3]. Among these complications, IPN is particularly
44 critical, occurring in 30%-40% of patients with severe AP and associated with a mortality rate
45 of 10%-20% [4, 5]. Therefore, early identification of high-risk IPN patients is crucial for timely
46 intervention and improved prognosis.

47 In the clinical management of AP, various scoring systems and biomarkers have been de-
48 veloped to predict disease severity. Traditional scoring systems, such as the Ranson score [6],

49 Bedside Index for Severity in Acute Pancreatitis (BISAP) [7], Modified CT Severity Index
50 (MCTSI) [8], and Acute Physiology and Chronic Health Evaluation II (APACHE II) [9], in-
51 tegrate multiple clinical indicators to provide valuable assessment tools for clinicians. Addi-
52 tionally, simplified scoring systems like the Pancreatitis Activity Scoring System (PASS) [10]
53 and the Chinese Severe Acute Pancreatitis Scoring System (CSSS) [11] have been introduced to
54 enhance the convenience and timeliness of early prediction. However, these systems exhibit sig-
55 nificant limitations in accurately predicting IPN, a severe complication of AP. They often lack
56 precision in identifying high-risk patients and fail to account for the unique pathophysiology of
57 IPN, resulting in suboptimal sensitivity and specificity. To address these limitations, machine
58 learning (ML) and deep learning (DL), as branches of artificial intelligence, have emerged as
59 powerful tools for disease prediction.

60 ML and DL models can integrate diverse clinical variables and improve predictive accuracy
61 compared to traditional methods [12]. In one study, logistic regression (LR), Fully connected
62 neural network (FCNN), and Extreme Gradient Boosting (XGBoost, XGB) were used to predict
63 the risk of severe acute pancreatitis (SAP), demonstrating that ML models outperformed tra-
64 ditional scoring systems in risk stratification [13]. Recent advances show that ML with feature
65 optimization and explainable AI (XAI) can achieve both high performance and transparency
66 across clinical domains. In dengue diagnosis, Recursive Feature Elimination with Shapley Ad-
67 ditive Explanations (SHAP) improved Random Forest (RF) performance while identifying key
68 predictors [21]. Similarly, in maternal health, Linear Discriminant Analysis (LDA)-optimized
69 features with ensemble learning achieved high accuracy, with SHAP clarifying decision logic
70 [22].

71 Evidence from hepatology further supports this approach. Studies using Deep Feature Syn-
72 thesis with tree ensembles achieved strong performance and provided interpretability through
73 SHAP analysis [23]. Other work combined tree-based selection with ensemble models and
74 XAI, though noting limitations like dataset bias and lack of external validation [25]. An Arti-
75 ficial Neural Network (ANN)-based framework optimized by LDA significantly improved accu-
76 racy while systematically integrating SHAP, Local Interpretable Model-agnostic Explanations
77 (LIME), and Individual Conditional Expectation (ICE), emphasizing the need for rigorous val-

78 idation [24].

79 Although ML has been increasingly applied in predicting disease severity, previous studies
80 on AP mainly focused on overall severity rather than IPN, and often lacked external validation
81 or model interpretability. To address these gaps, this study developed and externally validated
82 an interpretable ML framework for individualized IPN prediction. Multiple algorithms, includ-
83 ing linear, tree-based, and neural network models, were systematically compared under nested
84 cross-validation and temporal validation. Embedded feature selection methods (L1, Elastic Net,
85 and tree-based importance) and SHAP interpretation were integrated to enhance model robust-
86 ness and clinical explainability.

87 This approach combines methodological rigor and interpretability, distinguishing our work
88 from previous studies and providing a reliable tool for early IPN risk assessment.

89 **Method**

90 **Data Collection and Processing**

91 This was a single-center retrospective study conducted at Xuanwu Hospital, Capital Medical
92 University, including 728 patients with AP admitted between December 2017 and December
93 2023. Patients included in this study were 18 years of age or older and were diagnosed with
94 a first episode of AP according to the 2012 Atlanta Classification criteria. The classification
95 criteria for AP refer to the *Revised Atlanta Classification and Definitions by International Con-*
96 *sensus* [3]. Exclusion criteria included: (1) lack of early disease information, no abdominal
97 computed tomography (CT) scan within one week of onset, or no serum marker tests within
98 24 hours of onset; (2) patients with recurrent AP; (3) patients with chronic pancreatitis, pan-
99 creatitis related to trauma or pregnancy; (4) patients with cancer; (5) patients under 18 years
100 of age. After applying these criteria, 432 patients were classified into the IPN group, and 296
101 into the non-IPN group. A flow diagram of the selection process is presented in Supplementary
102 Figure S3.

103 Most predictor variables were collected within 24 hours of admission. Additional vari-
104 ables—such as comorbidities, MCTSI, and organ failure—were obtained at predefined time

105 points but always before IPN diagnosis. Since IPN typically developed days after admission,
 106 all variables were available prior to outcome determination. Details are provided in Supplemen-
 107 tary Table S3.

108 The primary outcome was IPN, defined according to the 2012 Revised Atlanta Classifica-
 109 tion and supported by the latest international consensus guidelines [3, 26]. Diagnosis required
 110 evidence of pancreatic or peripancreatic necrosis with either (1) microbiological confirmation
 111 or (2) gas within the necrotic area on contrast-enhanced computed tomography. The diagnosis
 112 date was defined as the earliest time either criterion was met.

113 This study was approved by the Ethics Committee of Xuanwu Hospital, Capital Medical
 114 University (Approval No.: XA Lin Yan Shen [KS2025] 002-001). All analyses were performed
 115 using anonymized data. Due to the retrospective nature of the study and the use of de-identified
 116 records, the requirement for individual informed consent was waived. All procedures were
 117 conducted in accordance with the principles of the Declaration of Helsinki.

118 In this study, we preprocessed the raw data to ensure its suitability for training and prediction
 119 with ML and FCNN. The main steps of preprocessing include binarization, one-hot encoding,
 120 normalization, and data distribution analysis.

121 Missing values were imputed using the median (for continuous variables) or mode (for cat-
 122 egorical variables). The proportion of missing data for each variable is provided in Supple-
 123 mentary Table S4. To prevent data leakage, all preprocessing steps, including imputation and
 124 normalization, were fitted exclusively on the training folds during cross-validation and subse-
 125 quently applied to the corresponding validation folds and the external test set. To address class
 126 imbalance, inverse class weighting was incorporated into the loss function, and stratified sam-
 127 pling was employed during cross-validation to ensure representative data splits. Categorical
 128 variables were one-hot encoded prior to model training, expanding several variables into multi-
 129 ple dummy features and resulting in a maximum of 33 input dimensions. Feature normalization
 130 was performed according to the following formula:

$$data[i]_{\text{norm}} = \frac{data[i] - data_{\min}}{data_{\max} - data_{\min}} \quad (1)$$

131 where $data[i]$ is the i -th data point, and $data_{\min}$ and $data_{\max}$ are the minimum and maximum

132 values, respectively.

133 All data preprocessing and model development were conducted in Python (v3.10) on a
134 Windows-based environment using commonly adopted scientific computing and machine-learning
135 libraries, including NumPy, pandas, scikit-learn, XGBoost, LightGBM, and Matplotlib/Seaborn
136 for visualization; SHAP analyses were performed using the SHAP Python package. All code
137 was executed in isolated environments to ensure reproducibility, and the exact package versions
138 were recorded programmatically and are provided as supplementary material or are available
139 upon reasonable request. All stochastic procedures used a fixed random seed, which is docu-
140 mented in the analysis scripts and available upon reasonable request.

141 **Model Development and Validation**

142 Six predictive models were developed: gradient boosting machine (GBM), XGB, RF, support
143 vector machine (SVM), LR, and a FCNN. The data were temporally split, with samples from
144 2017 to 2021 used for model development, and 2022 to 2023 held out as an independent external
145 test set for temporal validation. Nested cross-validation was applied within the training period:
146 the outer 5-fold CV was used to estimate model performance, while the inner 3-fold CV was
147 used for hyperparameter tuning and feature selection.

148 To mitigate data leakage, embedded feature selection was integrated directly into the model
149 training process within a nested cross-validation framework applied exclusively to the training
150 set (2017–2021). Within this framework, hyperparameters were optimized and feature subsets
151 were determined independently inside each training fold.

152 The final model was subsequently constructed by applying the optimized training procedure
153 to the entire training set, and its generalizability was evaluated through a single assessment on
154 the held-out temporal test set (2022–2023).

155 All models were implemented using scikit-learn and XGBoost, with class weighting applied
156 to address class imbalance. Hyperparameters were tuned using 5-fold inner cross-validation,
157 optimizing the area under the receiver operating characteristic curve (ROC-AUC).

158 For linear models, embedded feature selection was performed via L1 or Elastic Net regu-
159 larization in combination with SelectFromModel, retaining features with non-zero coefficients

Table 1: Clinical characteristics of patients (n = 728)

	IPN (n=432)	Non-IPN (n=296)	P value
Epidemiology			
Age, median (min, max), y	51.03 (11-90)	52.14 (14-99)	0.374
BMI, median (min, max), kg/m ²	23.89 (15.63-41.52)	24.62 (15.05-39.33)	0.024
Gender, n (%)			0.62
Male	282 (65.3%)	199 (67.2%)	
Female	150 (34.7%)	97 (32.8%)	
Medical History, n (%)			
Cardiovascular	178 (41.2%)	120 (40.5%)	
Cerebrovascular	25 (5.8%)	16 (5.4%)	
Respiratory	12 (2.8%)	11 (3.7%)	
Nephropathy	16 (3.7%)	9 (3.0%)	
Diabetes	78 (18.1%)	54 (18.2%)	
Substance Abuse, n (%)			
Smoking History	119 (27.5%)	99 (33.4%)	
Alcohol History	130 (30.1%)	101 (34.1%)	
Laboratory, Median (min, max)			
APACHE II	12.98 (1-36)	11.84 (1-34)	0.012
MCTSI	8.49 (2-10)	6.49 (2-10)	0.805
WBC (10 ⁹ /L)	11.13 (2.07-33.15)	11.70 (2.8-54.01)	0.230
Hematocrit	31.07 (12.2-55.9)	37.49 (16.4-57.3)	0.436
Lymph (10 ⁹ /L)	1.19 (0.07-15.94)	1.10 (0.14-3.29)	0.136
PLT (10 ⁹ /L)	282.13 (20-800)	247.67 (33-1084)	<0.001
Neut (10 ⁹ /L)	9.27 (1.21-30.23)	9.95 (1.17-51.33)	0.132
International Normalized Ratio	1.27 (0.9-1.93)	1.16 (0.87-2.52)	0.002
Activated Partial Thromboplastin Time (s)	42.93 (0-98.6)	40.43 (22.7-101.1)	<0.001
Thrombin Time (s)	16.17 (0-124.7)	16.12 (0-44.5)	0.897
Prothrombin Time (s)	15.79 (12.2-94.2)	14.81 (12-27.5)	<0.001
Fibrinogen (g/L)	4.92 (0.6-13.71)	5.28 (1.12-15)	0.011
D-Dimer (μg/mL)	4.55 (0-19.53)	4.51 (0-19.11)	0.909
Creatinine (μmol/L)	89.63 (10-953)	89.23 (15-987)	0.961
C-Reactive Protein (mg/L)	136.59 (2-571)	157.76 (3-471)	0.004
White Blood Interleukin-6 (pg/ml)	194.01 (0-5000)	242.15 (0-3666)	0.129
Procalcitonin (ng/ml)	4.03 (0-100)	4.17 (0-83.9)	0.871
Etiology			
Biliary	240 (55.6%)	156 (52.7%)	
Hypertriglyceridemia	127 (29.4%)	91 (30.7%)	
Alcoholic	18 (4.2%)	17 (5.7%)	
Other	47 (10.9%)	32 (10.8%)	
Complications, n (%)			
Respiratory Failure	122 (28.2%)	54 (18.2%)	
Renal Failure	81 (18.8%)	24 (8.1%)	
Heart Failure	91 (21.1%)	24 (8.1%)	

Data are presented as median (minimum–maximum) or count (percentage). Continuous variables were compared using the Mann–Whitney U test, and categorical variables were compared using the chi-square or Fisher's exact test as appropriate. P values < 0.05 were considered statistically significant. APACHE II, Acute Physiology and Chronic Health Evaluation II; MCTSI, Modified Computed Tomography Severity Index; WBC, White Blood Cell Count; Lymph, Lymphocyte Count; PLT, Platelet Count; Neut, Neutrophil Count.

160 within each training fold to prevent information leakage. The FCNN employed L1 or Elastic
 161 Net regularization at the input layer to encourage sparsity, depending on the configuration.

162 For tree-based models (RF, GBM, and XGBoost), no explicit feature elimination was con-
 163 ducted; instead, features were implicitly weighted according to split gain or impurity reduction
 164 during training. Importantly, all feature selection or weighting processes were confined strictly
 165 to the training folds.

166 The FCNN architecture consisted of one or two hidden layers containing 64 to 128 neurons
 167 per layer, with activation functions chosen from ReLU and Tanh. We performed hyperparame-
 168 ter optimization through a grid search coupled with five-fold cross-validation. The search em-
 169 ployed the Adam optimizer with learning rates ranging from 0.0005 to 0.001 and incorporated
 170 early stopping (patience=10). A summary of the hyperparameter search space and the resulting
 171 optimal configurations are provided in Table 2 and Supplementary Table S1, respectively.

Table 2: Architecture and hyperparameter search space of the FCNN.

Parameter	Range / Setting	Description
Input dimension	12, 15, 20, 31	Based on feature configuration
Hidden layers	(64,), (64, 32)	One or two layers
Activation function	ReLU, Tanh	Nonlinear transformation
Optimizer	Adam	Adaptive gradient optimization
Learning rate	0.0005, 0.001	Grid search range
Regularization (α)	0.001, 0.01, 0.1	L2 penalty
Max iterations	3000	Early stopping enabled
Patience	10 epochs	Validation monitored

The FCNN was tuned using a grid search within nested cross-validation, where hyperparameters were optimized based on ROC-AUC performance in the inner folds. Early stopping with a patience of 10 epochs was applied to prevent overfitting. The learning rate and regularization parameter (α) were selected according to validation performance. ReLU and Tanh activations were tested to balance convergence speed and non-linearity.

172 Model performance was primarily assessed using ROC-AUC. Secondary metrics included
 173 accuracy, precision, recall, F1-score, and Brier score. Ninety-five percent confidence intervals
 174 (95% CIs) for ROC-AUC, accuracy, and F1-score were estimated via 1000 bootstrap iterations.
 175 Calibration and precision–recall (PR) curves were plotted for visual assessment. The decision
 176 threshold for binary classification was determined by Youden’s Index, which maximizes the
 177 sum of sensitivity and specificity on the ROC curve. Specifically, the threshold was selected at
 178 the point where the true positive rate and the true negative rate were maximized. Pairwise ROC-

179 AUC differences between models were evaluated using a bootstrap-based DeLong-like test and
180 visualized as heatmaps. The best-performing model was further interpreted using SHAP values
181 to quantify feature contributions.

182 **Result**

183 **Model Performance and Feature Selection**

184 All models were trained using embedded feature selection within a nested cross-validation
185 framework, resulting in a spectrum of model complexities defined by the retention of 5 to 33
186 features. Tree-based ensembles—RF, GBM, and XGB—achieved the strongest external per-
187 formance, with ROC-AUC values approaching 0.76. The FCNN with elastic-net regularization
188 followed closely, achieving an ROC-AUC of 0.74, while linear models and the SVM delivered
189 slightly lower results, with ROC-AUCs in the range of 0.68 to 0.73. The robust generalizabil-
190 ity of all models was further confirmed by internal cross-validation, which yielded consistently
191 high mean ROC-AUCs between 0.81 and 0.82.

192 As a conventional clinical benchmark, the APACHE II score was also evaluated indepen-
193 dently on the external validation set. Its predictive performance was notably lower than that of
194 most ML models, with a ROC-AUC of 0.596 and a PR-AUC of 0.695, highlighting the potential
195 added value of multivariable ML approaches for IPN prediction.

196 Figure 1 displays the Receiver Operating Characteristic (ROC) and PR curves across all ML
197 models, while the PR curve of the APACHE II score is provided separately in the Supplementary
198 Materials. Comparative ROC-AUC and PR-AUC results are summarized in Figure 2. The
199 primary performance indicators, including ROC-AUC, accuracy, and F1-score, are reported
200 in Table 3, while additional metrics such as precision, recall, Brier score, and classification
201 thresholds are summarized in Supplementary Table S2.

202 **Statistical and Calibration Analyses**

203 Supplementary Figure S1 shows the results of statistical comparison and calibration analysis.
204 Subplots (a)–(c) present pairwise ROC-AUC comparisons using DeLong tests under different

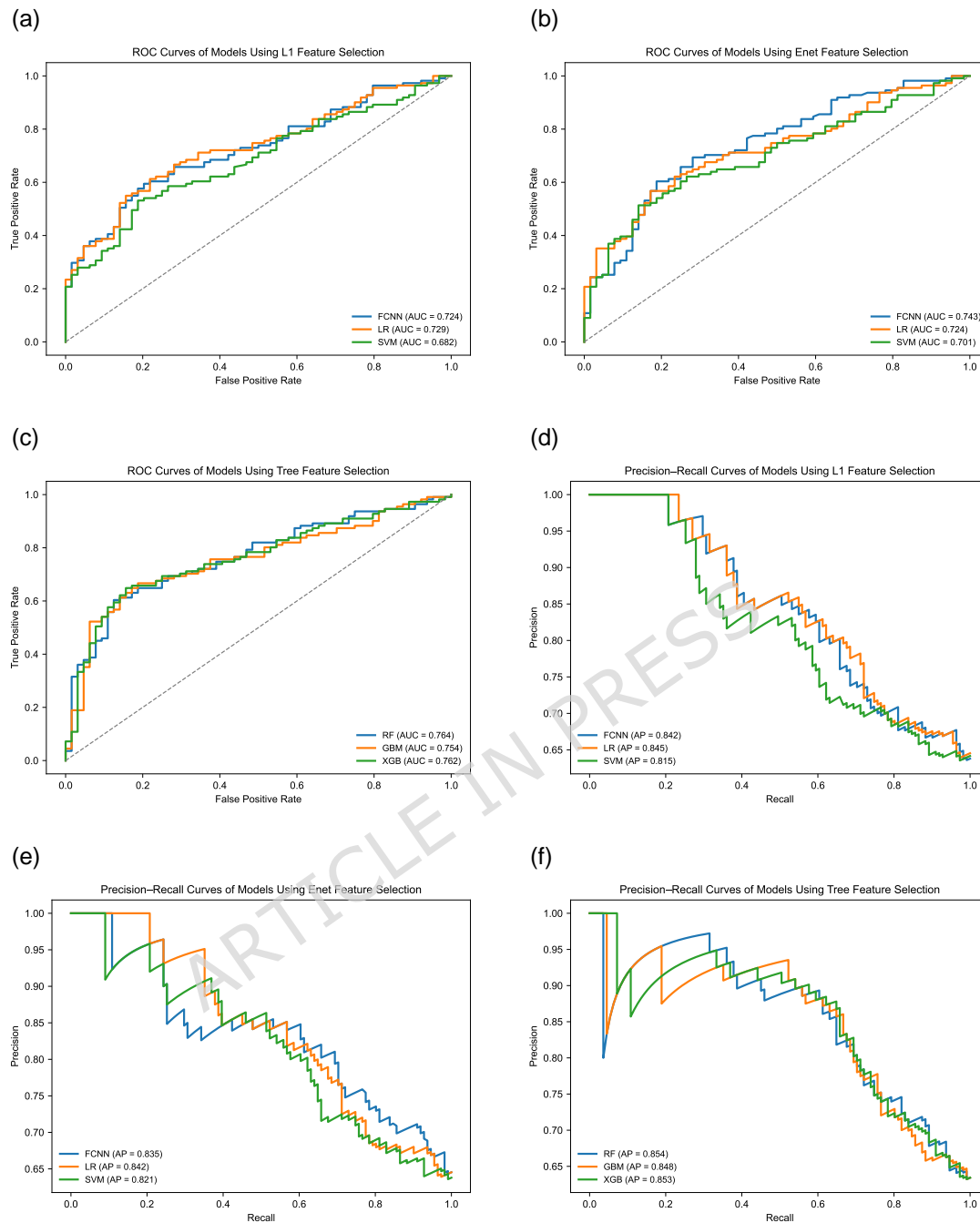


Figure 1: Model performance comparison by ROC-AUC and PR-AUC in the external validation cohort. Each panel illustrates the Receiver ROC and PR curves for models developed under three feature selection strategies: L1 regularization, Elastic Net, and tree-based embedded importance. Shaded areas represent 95% confidence intervals derived from 1000 bootstrap resamples. The decision threshold (marked by solid dots) corresponds to the Youden's Index, maximizing the sum of sensitivity and specificity. AUC values were compared across models using the bootstrap-based DeLong test.

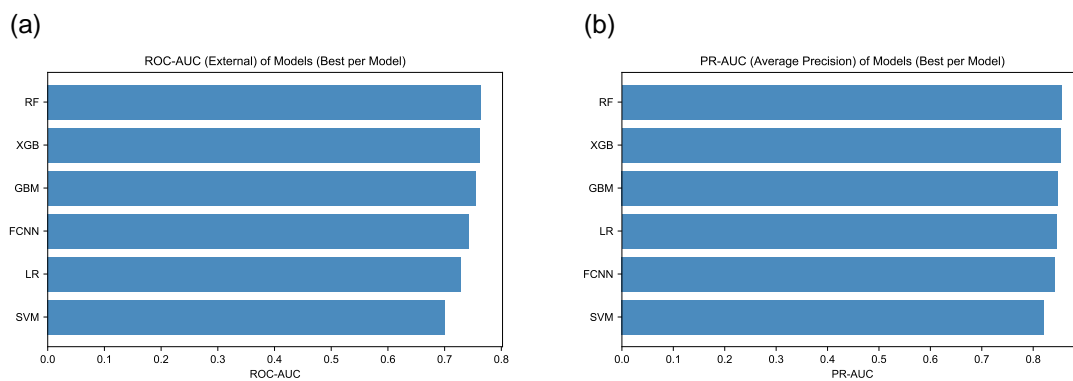


Figure 2: **ROC-AUC and PR-AUC performance across models.** (a) ROC-AUC values of each model on the external validation set. (b) PR-AUC (average precision) comparison across models.

Table 3: Performance metrics of all models on the external validation set.

Model	Features	Method	ROC-AUC	ROC-AUC (95% CI)	Accuracy	F1-score	PR-AUC
FCNN	9	embedded_l1	0.724	(0.653, 0.790)	0.669	0.695	0.842
LR	9	embedded_l1	0.729	(0.657, 0.797)	0.674	0.705	0.845
SVM	33	embedded_l1	0.682	(0.605, 0.755)	0.634	0.648	0.815
FCNN	5	embedded_enet	0.743	(0.664, 0.814)	0.680	0.705	0.835
LR	10	embedded_enet	0.724	(0.655, 0.794)	0.663	0.681	0.842
SVM	17	embedded_enet	0.701	(0.627, 0.775)	0.640	0.644	0.821
RF	33	embedded_tree	0.764	(0.696, 0.830)	0.703	0.720	0.854
GBM	33	embedded_tree	0.754	(0.682, 0.823)	0.720	0.751	0.848
XGB	33	embedded_tree	0.762	(0.687, 0.831)	0.720	0.746	0.853
APACHE II	N/A	scoring system	0.596	(0.507, 0.682)	0.657	0.752	0.695

205 feature selection methods. Tree-based models (RF, GBM, XGB) showed similar ROC-AUCs.
 206 In contrast, differences were observed between tree-based and linear or neural network models.
 207 Subplots (d)–(f) show calibration curves for the corresponding models. Tree-based models
 208 demonstrated better calibration, with predicted probabilities closer to the ideal line, compared
 209 to linear and neural network models.

210 Feature Interpretation of the Best Model

211 RF was selected due to its superior external ROC-AUC and well-calibrated performance. SHAP
 212 analysis revealed that higher levels of Fibrinogen, Prothrombin Time, APACHE II score, INR,
 213 and inflammatory markers such as IL-6 and C-reactive protein (CRP) were associated with an
 214 increased risk of IPN, as shown in Figure 3. Conversely, higher lymphocyte counts and lower
 creatinine or MCTSI values were linked to reduced IPN risk, indicating their protective roles.

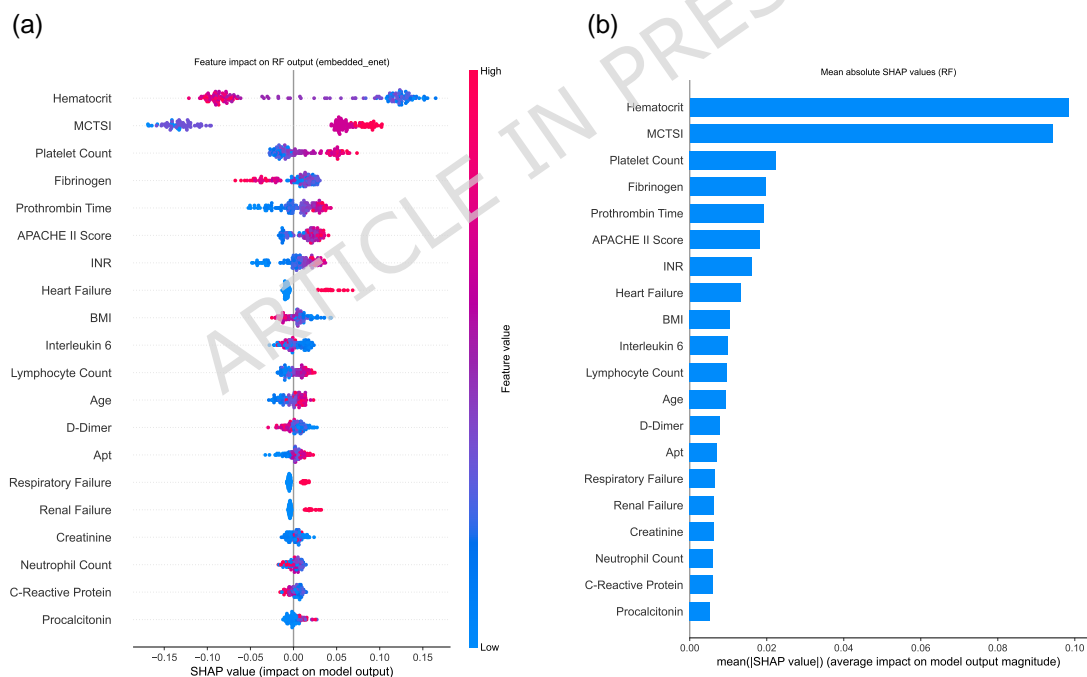


Figure 3: **SHAP-based feature interpretation of the best-performing RF model.** (a) Summary plot showing individual feature contributions to IPN prediction, where red indicates higher and blue lower feature values. (b) Bar plot ranking features by their mean absolute SHAP values, indicating overall feature importance.

216 Discussion

217 In this study, we systematically evaluated multiple ML models for predicting IPN using a ro-
218 bust framework that integrated embedded feature selection and temporal external validation.
219 A key finding was the superior and robust performance of tree-based ensemble methods, par-
220 ticularly RF. When trained on features selected by a tree-based selector, RF achieved optimal
221 performance on the independent test set, exhibiting excellent discrimination (AUC = 0.764)
222 and outstanding precision in identifying IPN cases (Average Precision = 0.854). Comparable
223 performance was observed for other tree-based models, including XGBoost and GBM, which
224 were evaluated under the same feature-selection paradigm. In contrast, models relying on L1- or
225 Elastic Net–based feature selection, such as LR, SVM, and FCNN, showed more moderate per-
226 formance. Collectively, these findings highlight the importance of aligning model architecture
227 with feature selection strategy and suggest that tree-based approaches are particularly effective
228 in capturing the complex, non-linear relationships underlying IPN development.

229 Earlier prediction models for AP severity and IPN mainly relied on multivariable LR using
230 clinician-selected variables [13, 14]. Although these approaches are transparent and easy to
231 interpret, they often fail to capture nonlinear and high-dimensional interactions in clinical data,
232 which limits their discriminative power and external generalizability.

233 More recent research has shifted toward machine and deep learning approaches that inte-
234 grate embedded feature selection, ensemble methods, and XAI to enhance predictive accuracy
235 and interpretability [16, 15]. Clinical studies across various diseases have demonstrated that
236 combining optimized feature selection with XAI techniques—such as SHAP and LIME—can
237 yield models that are both accurate and transparent, supporting more trustworthy decision-
238 making in medical settings [21, 23, 24].

239 Building on these methodological advances, our study demonstrates the added value of an
240 integrated ML framework for IPN prediction in acute pancreatitis. Among all evaluated models,
241 RF achieved the most favorable external performance, with superior calibration compared to
242 linear and neural network–based approaches. To contextualize the added value of ML, we
243 benchmarked model performance against the conventional APACHE II score. On the temporal
244 test set, APACHE II showed only moderate discrimination (ROC-AUC = 0.596, PR-AUC =

245 0.695), whereas most ML models achieved ROC-AUC values exceeding 0.70. These results
246 underscore the potential of multivariable learning-based approaches to enhance early IPN risk
247 stratification beyond traditional clinical scoring systems.

248 Beyond discrimination, calibration is essential for clinical deployment because it determines
249 whether predicted risks correspond to observed outcomes and supports threshold selection tai-
250 lored to clinical priorities. All tree-based ensembles produced clinically reliable probability
251 estimates, with the RF showing the lowest Brier score of 0.194, as reported in Supplementary
252 Table S2. At the standard threshold, RF also provided high precision (0.893) with moder-
253 ate recall (0.604), and decision curve analysis further supported its potential clinical utility.
254 In practice, excessive false-positive predictions may lead to unnecessary imaging or antibiotic
255 exposure, whereas false-negative predictions could delay timely intervention. Therefore, well-
256 calibrated risk estimates and flexible threshold adjustment are critical to balancing sensitivity
257 and specificity in real-world use.

258 Statistical comparisons via pairwise DeLong tests showed no statistically significant differ-
259 ences in discrimination among models within the same feature-selection group, as illustrated
260 in Supplementary Fig. S1(a)–(c), with all $p > 0.05$. Nevertheless, RF maintained a modest
261 numerical advantage over other tree-based ensembles and was further characterized by consis-
262 tently favorable calibration. These results support selecting RF as the most clinically robust
263 predictor, based on its practical consistency, reliable probability estimates, and interpretability,
264 rather than relying solely on marginal statistical superiority in AUC.

265 To further enhance model interpretability, SHAP analysis was conducted on the best-performing
266 RF model to quantify each feature's contribution to IPN prediction [17, 18]. The analysis re-
267 vealed that systemic inflammatory, coagulation, and organ dysfunction markers were the dom-
268 inant determinants of IPN, as illustrated in Figure 3. Elevated Fibrinogen, Prothrombin Time,
269 INR, CRP, Procalcitonin, Interleukin-6, and D-dimer were associated with increased predicted
270 risk, reflecting the hyperinflammatory and procoagulant state characteristic of severe acute pan-
271 creatitis. Conversely, higher Lymphocyte count, and Hematocrit, along with lower Creatinine
272 and MCTSI scores, were linked to reduced risk, suggesting preserved immune competence, ad-
273 equate perfusion, and less extensive parenchymal necrosis. Established severity indices such

274 as APACHE II and MCTSI also ranked among the most influential predictors, demonstrating
275 strong concordance between model-derived and clinical understanding.

276 From a clinical standpoint, the SHAP-derived pattern closely mirrors the pathophysiologi-
277 cal cascade leading from sterile necrosis to infection. Elevated coagulation and inflammatory
278 markers indicate endothelial injury and microthrombosis, which compromise pancreatic perfu-
279 sion and create a favorable environment for bacterial translocation. Meanwhile, lower Hemat-
280 ocrit values reflect hypovolemia and catabolic stress, often observed in patients who progress to
281 IPN. Representative features such as MCTSI, Hematocrit, and D-dimer thus bridge radiologic
282 severity, circulatory disturbance, and coagulation dysfunction—the three interdependent mech-
283 anisms underlying infected necrosis. This coherence between model insights and biological
284 mechanisms reinforces both the interpretability and credibility of the predictive framework.

285 The robustness and generalizability of the models were supported by consistent internal and
286 external performance. During cross-validation, all models achieved mean AUC values between
287 0.81 and 0.82, with RF and XGB showing particularly stable performance at 0.816 ± 0.020 and
288 0.817 ± 0.011 , respectively. Detailed cross-validation outcomes are provided in Supplementary
289 Table S2. External validation revealed only a minor decline in discrimination, demonstrating
290 good temporal transportability and model stability. These findings indicate that the proposed
291 framework can generalize reliably to new patient populations, despite inherent heterogeneity
292 in disease presentation and clinical practice. However, as a single-center retrospective study,
293 this work remains limited by potential selection bias, residual confounding, and temporal drift.
294 Future multicenter and prospective validation will be essential to confirm its reproducibility and
295 ensure real-world applicability before clinical deployment.

296 Clinically, the proposed ML framework has the potential to support early risk stratification
297 of IPN, enabling timely identification of high-risk patients and optimization of clinical decision-
298 making. Such predictive tools could guide escalation of care, such as earlier transfer to intensive
299 monitoring, targeted antibiotic therapy, or closer imaging surveillance in patients predicted to
300 develop infection. By integrating routinely available laboratory and clinical parameters, the
301 model is both interpretable and feasible for real-time application in tertiary and emergency care
302 settings.

303 This study has several limitations that should be acknowledged. First, this study included
304 a moderate sample size ($n = 728$) with 432 IPN events. As Xuanwu Hospital is a large ter-
305 tiary referral center, the case mix may be enriched for more severe presentations, which could
306 contribute to a relatively higher IPN proportion compared with primary or community settings.
307 Second, although we implemented nested cross-validation, embedded feature selection, and an
308 independent temporal validation cohort, the possibility of model optimism due to center-specific
309 practice patterns cannot be completely excluded. Importantly, the consistent performance across
310 internal cross-validation and temporal validation suggests that overfitting was mitigated to a
311 substantial extent. Third, while temporal external validation was performed, validation across
312 different institutions is still needed to confirm transportability under varying laboratory assays,
313 clinical workflows, and patient populations. Prospective multicenter studies and, if necessary,
314 model recalibration are warranted before broader clinical implementation.

315 Future work should prioritize validation in multicenter cohorts spanning different care set-
316 tings, including both tertiary referral and non-referral hospitals, to further assess transportability
317 and calibration stability across institutions. Model recalibration and, where appropriate, transfer
318 learning approaches may help adapt predictions to evolving clinical standards, laboratory as-
319 says, and population characteristics over time. In addition, integrating multimodal data—such
320 as radiological features and longitudinal trajectories of inflammatory markers—may further en-
321 hance predictive performance and provide deeper mechanistic insight. Finally, future studies
322 should extend evaluation beyond discrimination and calibration to assess real-world clinical
323 impact, including outcome-oriented and cost-sensitive analyses, to determine whether model-
324 assisted decision-making improves patient outcomes compared with standard practice.

325 **Conclusion**

326 This study developed and evaluated an interpretable ML framework for early prediction of IPN
327 in patients with acute pancreatitis, with IPN defined according to established diagnostic crite-
328 ria. Predictors were restricted to routinely available variables collected early after admission and
329 prior to IPN diagnosis to minimize information leakage, and embedded feature selection was
330 integrated within a nested cross-validation framework. Using a temporal hold-out cohort for val-

331 idation, the random forest model achieved the best overall performance, with an external AUC
332 of 0.764, balanced precision and recall, and the lowest Brier score, indicating reliable probabilit-
333 ity calibration. SHAP analysis highlighted clinically plausible predictors, including fibrinogen,
334 APACHE II score, D-dimer, IL-6, and CRP as risk-enhancing factors, whereas lymphocyte
335 count and hematocrit were associated with reduced risk. These results support the potential
336 of transparent, tree-based models for early, data-driven risk stratification in acute pancreatitis.
337 Future work should prioritize prospective multicenter validation, calibration assessment and re-
338 calibration across institutions, and integration of multimodal data to enhance generalizability
339 and real-world applicability.

340 **Abbreviations**

AP	Acute pancreatitis
IPN	Infected pancreatic necrosis
BISAP	Bedside Index for Severity in Acute Pancreatitis
MCTSI	Modified Computed Tomography Severity Index
APACHE II	Acute Physiology and Chronic Health Evaluation II
PASS	Pancreatitis Activity Scoring System
CSSS	Chinese Simple Severity Score
ML	Machine learning
DL	Deep learning
LR	Logistic regression
FCNN	Fully connected neural network
XGB	Extreme gradient boosting
SAP	Severe acute pancreatitis
XAI	Explainable artificial intelligence
341 SHAP	Shapley additive explanations
RF	Random forest
LDA	Linear discriminant analysis
ANN	Artificial neural network
LIME	Local Interpretable Model-agnostic Explanations
ICE	Individual Conditional Expectation
CT	Computed tomography
GBM	Gradient boosting machine
SVM	Support vector machine
ROC-AUC	Area under the receiver operating characteristic curve
PR	Precision–recall
ROC	Receiver operating characteristic
Neut	Neutrophil count
INR	International normalized ratio

342 **Declarations**

343 **Ethics approval and consent to participate**

344 This retrospective study was approved by the Ethics Committee of Xuanwu Hospital, Capital
345 Medical University (Approval No.: XA Lin Yan Shen [KS2025] 002-001). The study used
346 pre-existing, fully anonymized clinical data; therefore, the requirement for individual informed
347 consent was waived by the ethics committee.

348 **Consent for publication**

349 Not applicable.

350 **Data Availability Statement**

351 The datasets generated and/or analyzed during the current study are not publicly available due
352 to institutional data use policies but are available from the corresponding author on reasonable
353 request.

354 **Competing interests**

355 The authors declare that they have no competing interests

356 **Funding**

357 This work was financially supported by two grants: the Hebei Natural Science Foundation
358 (Grant No. H2024112019) and the S&T Program of Xiongan New Area (Grant No. XA202401102001K).

359 **Authors' contributions**

360 Xin Li was responsible for code implementation, data preprocessing, model optimization, manuscript
361 writing, and submission. Yixuan Ding contributed to clinical data collection, assisted with
362 data preprocessing, and participated in manuscript revision and polishing. Bohan Huang par-
363 ticipated in data collection, developed inclusion criteria, and contributed to data processing.
364 Yunheng Shen assisted with early-stage code development and model debugging. Hairong Lv
365 supported model tuning and performance optimization. Feng Cao provided clinical supervision,
366 contributed to the study design, and revised the manuscript. Tong Yu polished the manuscript
367 and contributed to language refinement. Fei Li, Xiaolu Fei, and Jia Li served as corresponding
368 authors, supervised the overall project, and provided critical revision of the manuscript.

369 **Acknowledgements**

370 No acknowledgements.

Author details

¹Department of Biomedical Engineering, Capital Medical University, No. 10, Youanmen Wai Xitoutiao, Beijing 100069, China.

²Tsinghua University, No. 1, Tsinghua Garden, Beijing 100084, China.

³Xuanwu Hospital of Capital Medical University, No. 45, Changchun Street, Beijing 100053, China.

⁴Xiongan Xuanwu Hospital, No. 1, Xili Road, Xufa District, Xiongan New Area, Hebei 070001, China.

References

[1] Peery, A. F. et al. Burden of gastrointestinal disease in the United States: 2012 update. *Gastroenterology* **143**, 1179–1187.e3 (2012). <https://doi.org/10.1053/j.gastro.2012.08.002>

[2] Garg, S. K. et al. Incidence, admission rates, and economic burden of adult emergency visits for chronic pancreatitis: Data from the National Emergency Department Sample, 2005 to 2012. *J. Clin. Gastroenterol.* **53**, e328–e333 (2019). <https://doi.org/10.1097/MCG.0000000000001096>

[3] Banks, P. A. et al. Classification of acute pancreatitis—2012: Revision of the Atlanta classification and definitions by international consensus. *Gut* **62**, 102–111 (2013). <https://doi.org/10.1136/gutjnl-2012-302779>

[4] Shah, J., Fernandez Y Viesca, M., Jagodzinski, R. and Arvanitakis, M. Infected pancreatic necrosis—current trends in management. *Indian J. Gastroenterol.* **43**, 578–591 (2024). <https://doi.org/10.1007/s12664-023-01506-w>

[5] de-Madaria, E. and Buxbaum, J. L. Advances in the management of acute pancreatitis. *Nat. Rev. Gastroenterol. Hepatol.* **20**, 691–692 (2023). <https://doi.org/10.1038/s41575-023-00808-w>

- 396 [6] Kumar, A. H. and Griwan, M. S. A comparison of APACHE II, BISAP, Ran-
397 son's score and modified CTSI in predicting the severity of acute pancreatitis based
398 on the 2012 revised Atlanta classification. *Gastroenterol. Rep.* **6**, 127–131 (2018).
399 <https://doi.org/10.1093/gastro/gox029>
- 400 [7] Gao, W., Yang, H. X. and Ma, C. E. The value of BISAP score for predicting mortality
401 and severity in acute pancreatitis: A systematic review and meta-analysis. *PLoS One* **10**,
402 e0130412 (2015). <https://doi.org/10.1371/journal.pone.0130412>
- 403 [8] Alberti, P. et al. Evaluation of the modified computed tomography severity in-
404 dex (MCTSI) and computed tomography severity index (CTSI) in predicting sever-
405 ity and clinical outcomes in acute pancreatitis. *J. Dig. Dis.* **22**, 41–48 (2021).
406 <https://doi.org/10.1111/1751-2980.12961>
- 407 [9] Knaus, W. A., Draper, E. A., Wagner, D. P. and Zimmerman, J. E. APACHE II: a severity
408 of disease classification system. *Crit. Care Med.* **13**, 818–829 (1985).
- 409 [10] Mao, W. et al. Prediction of infected pancreatic necrosis in acute necrotizing pancreatitis
410 by the modified pancreatitis activity scoring system. *United Eur. Gastroenterol. J.* **11**,
411 69–78 (2023). <https://doi.org/10.1002/ueg2.12353>
- 412 [11] Wang, L. et al. A simple new scoring system for predicting the mortality of severe acute
413 pancreatitis: a retrospective clinical study. *Medicine (Baltimore)* **99**, e20646 (2020).
414 <https://doi.org/10.1097/MD.00000000000020646>
- 415 [12] Sadr, H., Nazari, M., Khodaverdian, Z. et al. Unveiling the potential of artificial in-
416 telligence in revolutionizing disease diagnosis and prediction: a comprehensive review
417 of machine learning and deep learning approaches. *Eur. J. Med. Res.* **30**, 418 (2025).
418 <https://doi.org/10.1186/s40001-025-02680-7>
- 419 [13] Thapa, R. et al. Early prediction of severe acute pancreatitis using machine learning.
420 *Pancreatology* **22**, 43–50 (2022). <https://doi.org/10.1016/j.pan.2021.10.003>
- 421 [14] Wiese, M. L. et al. Identification of early predictors for infected necrosis in acute pancre-
422 atitis. *BMC Gastroenterol.* **22**, 405 (2022). <https://doi.org/10.1186/s12876-022-02490-9>

- 423 [15] Zhang, H. et al. Tree-based ensemble machine learning models in the prediction of acute
424 respiratory distress syndrome following cardiac surgery: a multicenter cohort study. *J.*
425 *Transl. Med.* **22**, 772 (2024). <https://doi.org/10.1186/s12967-024-05395-1>
- 426 [16] Muhammad, D. and Bendeche, M. Unveiling the black box: a systematic review of
427 explainable artificial intelligence in medical image analysis. *Comput. Struct. Biotechnol.*
428 *J.* **24**, 542–560 (2024). <https://doi.org/10.1016/j.csbj.2024.08.005>
- 429 [17] Shakeri, E. et al. Explaining eye diseases detected by machine learning using SHAP: a
430 case study of diabetic retinopathy and choroidal nevus. *SN Comput. Sci.* **4**, 433 (2023).
431 <https://doi.org/10.1007/s42979-023-01859-1>
- 432 [18] Rao, S., Mehta, S., Kulkarni, S., Dalvi, H., Katre, N. and Narvekar, M. A
433 study of LIME and SHAP model explainers for autonomous disease predictions.
434 *Proc. IEEE Bombay Sect. Signature Conf. (IBSSC)*, Mumbai, India, 1–6 (2022).
435 <https://doi.org/10.1109/IBSSC56953.2022.10037324>
- 436 [19] Sahu, B., Abbey, P., Anand, R., Kumar, A., Tomer, S. and Malik, E. Severity assessment
437 of acute pancreatitis using CT severity index and modified CT severity index: correlation
438 with clinical outcomes and severity grading as per the Revised Atlanta classification.
439 *Indian J. Radiol. Imaging* **27**, 152–160 (2017). https://doi.org/10.4103/ijri.IJRI_300_16
- 440 [20] Wan, J. et al. Serum D-dimer levels at admission for prediction of outcomes in acute pan-
441 creatitis. *BMC Gastroenterol.* **19**, 67 (2019). <https://doi.org/10.1186/s12876-019-0989-x>
- 442 [21] Das, K. et al. Optimized feature-driven dengue diagnosis using explain-
443 able machine learning approaches. *Proc. Int. Conf. Quantum Photonics, Ar-*
444 *tificial Intelligence, and Networking (QPAIN)*, Rangpur, Bangladesh, **1–6**
445 (2025). <https://doi.org/10.1109/QPAIN66474.2025.11171726>
- 446 [22] Mamun, M., Hussain, M. I., Ali, M. S., Alam Chowdhury, M. S., Chowdhury, S. H.
447 and Hossain, M. M. An explainable ensemble learning framework with feature opti-
448 mization for accurate maternal health risk prediction. *Proc. Int. Conf. Quantum Photon-*

- 449 *ics, Artificial Intelligence, and Networking (QPAIN)*, Rangpur, Bangladesh, **1–6** (2025).
450 <https://doi.org/10.1109/QPAIN66474.2025.11172243>
- 451 [23] Chowdhury, S. H. et al. Hepatitis C detection from blood donor data using hybrid
452 deep feature synthesis and interpretable machine learning. *Proc. 2nd Int. Conf. Next-*
453 *Generation Computing, IoT and Machine Learning (NCIM)*, Gazipur, Bangladesh, **1–6**
454 (2025). <https://doi.org/10.1109/NCIM65934.2025.11160156>
- 455 [24] Chowdhury, S. H., Mamun, M., Shaikat, T. A., Hussain, M. I., Iqbal, S. and Hossain, M.
456 M. An ensemble approach for artificial neural network-based liver disease identification
457 from optimal features through hybrid modeling integrated with advanced explainable AI.
458 *Medinformatics* **2**, 107–119 (2025). <https://doi.org/10.47852/bonviewMEDIN52024744>
- 459 [25] Mamun, M., Chowdhury, S. H., Hossain, M. M., Khatun, M. R. and Iqbal, S.
460 Explainability-enhanced liver disease diagnosis technique using tree selection and stack-
461 ing ensemble-based random forest model. *Informatics and Health* **2**, 17–40 (2025).
462 <https://doi.org/10.1016/j.infoh.2025.01.001>
- 463 [26] Párnitzky, A. et al. International Association of Pancreatology Revised Guidelines on
464 Acute Pancreatitis 2025: Supported and Endorsed by the American Pancreatic Asso-
465 ciation, European Pancreatic Club, Indian Pancreas Club, and Japan Pancreas Society.
466 *Pancreatology* **25**(6), 770–814 (2025). <https://doi.org/10.1016/j.pan.2025.04.020>