



A machine learning-based prediction model for colorectal liver metastasis

Sisi Feng¹ · Manli Zhou¹ · Zixin Huang¹ · Xiaomin Xiao¹ · Baiyun Zhong^{1,2}

Received: 1 April 2025 / Accepted: 15 April 2025
© The Author(s) 2025

Abstract

Colorectal liver metastasis (CRLM) is a primary factor contributing to poor prognosis and metastasis in colorectal cancer (CRC) patients. This study aims to develop and validate a machine learning (ML)-based risk prediction model using conventional clinical data to forecast the occurrence of CRLM. This retrospective study analyzed the clinical data of 865 CRC patients between January 2018 and September 2024. Patients were categorized into non-CRLM and CRLM groups. The least absolute shrinkage and selection operator regression was employed to identify key clinical variables, and five ML algorithms were utilized to develop prediction models. The optimal model was selected based on performance metrics including the receiver operating characteristic curve, precision-recall curve, decision curve analysis, and calibration curve, which collectively evaluated both the predictive accuracy and clinical utility of the model. Among the five ML algorithms evaluated, Random forest demonstrated the best performance. Leveraging the Random forest algorithm, we developed the CRLM-Lab6 prediction model, which incorporates six features: LDH, CA199, ALT, CEA, TBIL, and AGR. This model exhibits robust predictive performance, achieving an area under the curve of 0.94, a sensitivity of 0.88, and a specificity of 0.93. To enhance its practical utility, the model has been integrated into an accessible web application. This study developed a novel risk prediction model by integrating ML algorithms with conventional laboratory test data to evaluate the likelihood of CRLM occurrence. The model demonstrates excellent predictive performance and has significant clinical application potential.

Keywords Colorectal cancer · Liver metastasis · Machine learning · Prediction model · Screening

Introduction

Colorectal cancer (CRC) is one of the most prevalent malignant tumors of the digestive tract. According to global statistics, CRC ranks third in incidence and second in mortality among all cancers, with both rates showing an upward trend annually [1]. Approximately 30–40% of CRC patients present with metastasis at the time of diagnosis [2, 3]. The liver is the most common site for distant metastasis in CRC

patients [4]. Within five years of diagnosis, approximately 12.8% of CRC patients develop colorectal liver metastasis (CRLM), which is a significant factor contributing to treatment failure and patient mortality [5–7]. Despite advancements in imaging techniques such as CT, MRI, and PET-CT, the diagnostic accuracy for smaller metastatic lesions remains limited [8, 9]. Moreover, conventional tumor markers like carcinoembryonic antigen (CEA) and carbohydrate antigen 199 (CA199) exhibit suboptimal sensitivity and specificity for detecting CRLM [10]. Therefore, exploring new methods and markers for early diagnosis of CRLM is crucial.

With the rapid advancement of information technology and the digital transformation in the medical field, the volume and diversity of medical data are expanding at an unprecedented rate. These data primarily originate from laboratory information systems (LIS) and electronic medical record (EMR) systems [11]. Numerous studies have investigated the predictive and prognostic significance of individual variables within this dataset for CRLM patients,

Sisi Feng and Manli Zhou have contributed equally to this work and should be considered as co-first authors.

✉ Baiyun Zhong
xycsuhn@163.com

¹ Department of Clinical Laboratory, Xiangya Hospital, Central South University, Changsha 410008, Hunan, China

² National Clinical Research Center for Geriatric Disorders, Xiangya Hospital, Central South University, Changsha 410008, Hunan, China

such as CEA, alanine aminotransferase (ALT), and bilirubin levels [12–14]. Notably, these variables often exhibit intricate interdependencies. Unraveling these relationships through various analytical strategies is crucial for enhancing the quality of medical services and optimizing the allocation of healthcare resources [15]. Machine learning (ML), a subset of artificial intelligence, holds significant potential in aiding the diagnosis of CRLM [16]. ML algorithms can efficiently process and analyze large-scale, high-dimensional medical datasets, uncovering hidden disease patterns and associations that provide robust support for early disease detection [17, 18]. For instance, Miller et al. [19] developed a Random forest model integrating metabolomics analysis with demographic data (age and gender), achieving a specificity of 0.94 in the CRLM cohort, albeit with limited sensitivity at 0.51.

In this study, we evaluate the predictive performance of multiple ML models to identify risk factors associated with CRLM and offer novel insights for its early clinical identification and diagnosis.

Patients and methods

Study population

This study was conducted at Xiangya Hospital, Central South University. We retrospectively analyzed 613 CRC patients from January 2018 to October 2022, comprising 373 non-CRLM patients and 240 CRLM patients. These patients were included in the discovery cohort for model construction. For the validation cohort, we prospectively collected data from 252 subjects between January 2023 and September 2024, including 160 non-CRLM patients, and 92 CRLM patients, to evaluate the predictive performance of the optimal model. The inclusion criteria were as follows: (1) histopathological confirmation of CRLM or non-CRLM; (2) no prior radiotherapy, chemotherapy, or hormone therapy before admission. Tumor staging was performed according to the tumor-node-metastasis (TNM) classification system [20]. The exclusion criteria were: (1) concurrent presence of other malignant tumors; (2) severe dysfunction of vital organs such as the heart and lung; (3) pregnancy or lactation.

Data collection

We collected the basic information of all first laboratory tests after admission from the LIS and EMR databases. The data included: (1) demographic characteristics such as age and gender; (2) Laboratory indicators, including white blood cell (WBC), hemoglobin (Hb), platelet count, red cell distribution width (RDW), mean corpuscular volume (MCV), mean platelet volume (MPV), albumin-to-globulin ratio

(AGR), total bilirubin (TBIL), direct bilirubin (DBIL), total bile acids (TBA), ALT, aspartate aminotransferase (AST), lactate dehydrogenase (LDH), creatinine, creatine kinase (CK), myoglobin (Mb), triglycerides (TG), low-density lipoprotein cholesterol (LDL-C), high-density lipoprotein cholesterol (HDL-C), CEA and CA199; (3) Pathological data, including tumor size, location, T stage, etc. Additionally, in this study, we employed the least absolute shrinkage and selection operator (LASSO) regression to screen variables for incorporation into the ML model. LASSO is a regression method that utilizes L1 regularization to introduce a penalty term, thereby enabling variable selection, reducing model complexity, and mitigating the risk of overfitting.

Data cleaning

(1) Data preprocessing: Exclude detection data with a missing rate exceeding 30%. Based on the characteristics of the data distribution, select appropriate statistical methods that accurately reflect the central tendency of variables (e.g., median, mean) to impute missing values; (2) Data normalization: Standardize and normalize the four critical elements—specimen type, test item name, test result unit, and test reference range, to mitigate the impact of dimensional differences and magnitude disparities among features. This enhances the convergence speed and predictive accuracy of the algorithm.

Model construction and evaluation

The workflow for constructing the model using ML algorithms is illustrated in Fig. 1. In this study, five distinct ML models were employed for training and validation: Logistic regression, Linear support vector classification (Linear SVC), Random forest, Decision tree, and Support vector machine (SVM). The dataset was randomly divided into a training set and a validation set at a ratio of 7:3. The training set was utilized to develop the ML models, while the validation set served to assess their predictive performance. Fivefold cross-validation was implemented to identify the optimal hyperparameters for each of the five ML models. The primary metric for evaluating model performance was the area under the receiver operating characteristic (ROC) curve (AUC). Additionally, we conducted an in-depth analysis of several other performance metrics, including sensitivity, specificity, accuracy, precision, F1 score, positive predictive value (PPV), and negative predictive value (NPV). To provide a more comprehensive evaluation of the models, we also performed precision-recall curve analysis, decision curve analysis, and calibration curve assessment to evaluate both the clinical utility and the consistency between predicted probabilities and actual outcomes. Model construction and visualization were conducted using the Deepwise

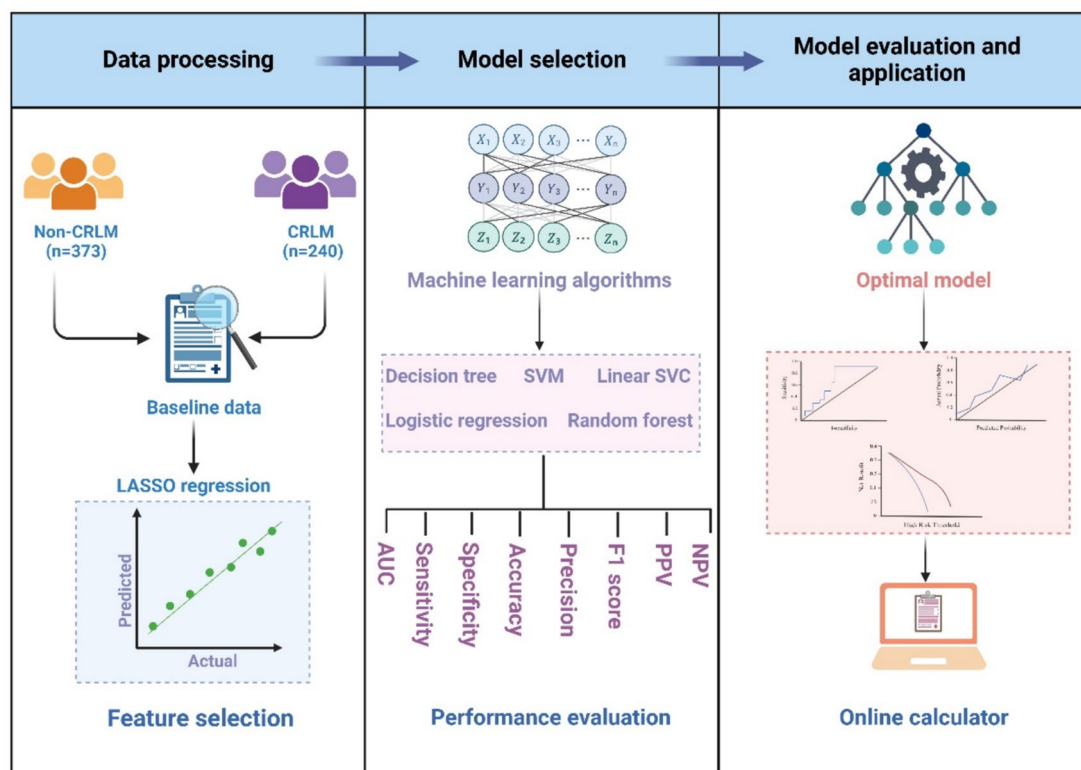


Fig. 1 Flowchart of ML for building the CRLM prediction model

and Beckman Coulter DxAI platform (<https://dxonline.deepwise.com/login>).

Statistical analysis

Statistical analysis was conducted using SPSS 23.0 and R 3.6.1. Continuous variables following a normal distribution were expressed as mean \pm standard deviation, whereas those not conforming to a normal distribution were reported as median (interquartile range). Categorical variables were summarized using percentages. Intergroup comparisons were performed using Student's t-test, Mann–Whitney U test, or chi-square test, as appropriate. A two-sided *P* value of less than 0.05 was considered statistically significant.

Results

Demographic baseline data

Table 1 summarizes the clinical characteristics of participants in both the discovery and validation cohorts. A total of 865 subjects were included in this study. The discovery cohort comprised 613 CRC patients, including 373

non-CRLM and 240 CRLM patients. The validation cohort consisted of 252 subjects, including 160 non-CRLM and 92 CRLM patients. In the discovery cohort, no statistically significant differences were observed between the non-CRLM and CRLM groups in terms of age, gender, WBC count, Hb, platelet, RDW, MPV, creatinine, CK, Mb, HDL-C, tumor size, location, and degree of differentiation ($P > 0.05$). However, CRLM patients exhibited more frequent abnormalities in liver function, characterized by elevated levels of TBIL, DBIL, TBA, ALT, and AST, as well as increased tumor markers CEA and CA199. Additionally, there were significant differences between the two groups in MCV, TG, AGR, LDH, LDL-C, T stage, and lymph node metastasis ($P < 0.05$). The data distribution between the discovery and validation cohorts was largely consistent.

Feature selection

To further enhance the interpretability and predictive performance of the model, we employed LASSO regression for variable selection. Using tenfold cross-validation, we identified nine features with non-zero coefficients: LDH, CA199, ALT, CEA, TBIL, AGR, TG, MCV, and lymph node metastasis (Fig. 2a, b).

Table 1 Baseline characteristics of CRLM and non-CRLM patients

Variables	Discovery cohort			Validation cohort		
	Non-CRLM	CRLM	<i>P</i>	Non-CRLM	CRLM	<i>P</i>
	N = 373	N = 240		N = 160	N = 92	
<i>Demographics</i>						
Age (Year)	61.0 (54.0–69.0)	60.0 (53.0–68.0)	0.241	64.0 (55.0–68.0)	64.0 (58.0–66.0)	0.331
Gender, Female/Male [n (%)]	162/211 (43.4/56.6)	91/149 (37.9/62.1)	0.176	78/82 (48.8/51.2)	38/54 (41.3/58.7)	0.254
<i>Laboratory routine testing</i>						
WBC (10 ⁹ /L)	5.7 (4.8–7.2)	6.1 (4.7–7.9)	0.125	5.4 (4.3–6.6)	5.6 (4.5–7.5)	0.301
Hb (g/L)	116.0 (95.0–130.0)	117.0 (103.8–131.0)	0.122	113.0 (90.0–130.3)	115.5 (101.0–127.0)	0.740
Platelet (10 ⁹ /L)	246.0 (191.0–298.0)	247.5 (184.8–316.0)	0.931	246.0 (187.0–298.0)	216.0 (173.0–274.0)	0.032
RDW (%)	14.5 (13.4–17.3)	14.7 (13.5–17.0)	0.547	14.2 (13.4–17.5)	14.7 (13.8–17.8)	0.115
MCV (fl)	88.1 (82.0–93.3)	90.0 (84.6–94.2)	0.007	88.3 (81.0–92.4)	93.7 (90.4–97.9)	<0.001
MPV (fl)	8.6 (7.9–9.4)	8.5 (7.9–9.5)	0.812	8.4 (7.8–9.1)	8.7 (8.1–9.8)	0.032
AGR	1.4 (1.3–1.6)	1.2 (1.1–1.4)	<0.001	1.4 ± 0.3	1.3 ± 0.2	0.021
TBIL (μmol/L)	9.5 (6.6–14.0)	13.5 (8.6–20.1)	<0.001	9.0 (6.6–12.5)	15.5 (10.6–19.7)	<0.001
DBIL (μmol/L)	2.7 (2.1–3.8)	3.9 (2.7–6.9)	<0.001	2.8 (2.1–3.7)	3.8 (3.0–5.0)	<0.001
TBA (μmol/L)	3.3 (1.9–5.8)	4.3 (2.1–7.5)	0.003	3.3 (2.1–5.4)	5.3 (2.9–8.8)	<0.001
ALT (U/L)	14.0 (10.1–23.0)	23.3 (14.8–37.5)	<0.001	14.5 (10.5–26.6)	30.5 (20.6–36.5)	<0.001
AST (U/L)	19.2 (15.5–25.6)	28.7 (20.4–44.2)	<0.001	19.7 (16.8–25.6)	24.6 (18.6–41.1)	<0.001
LDH (U/L)	171.0 (145.0–200.0)	270.5 (222.8–353.0)	<0.001	185.0 (155.9–212.5)	277.0 (195.0–388.3)	<0.001
Creatinine (μmol/L)	65.0 (54.0–79.0)	68.0 (56.2–78.6)	0.189	63.0 (51.0–77.0)	58.1 (46.0–72.1)	0.037
CK (U/L)	56.2 (38.7–75.6)	57.8 (41.0–80.7)	0.347	53.5 (36.9–69.7)	52.6 (37.6–78.5)	0.700
Mb (ug/L)	40.5 (29.2–52.4)	39.8 (30.1–50.0)	0.523	38.7 (30.1–49.4)	34.6 (26.5–43.8)	0.056
TG (mmol/L)	1.3 (0.9–1.8)	1.5 (1.1–2.2)	<0.001	1.4 (1.1–1.9)	2.3 (1.5–2.8)	<0.001
HDL-C (mmol/L)	1.1 (0.9–1.2)	1.1 (0.8–1.2)	0.177	1.0 (0.8–1.2)	0.9 (0.6–1.1)	0.051
LDL-C (mmol/L)	2.9 (2.4–3.3)	3.0 (2.5–3.7)	0.003	2.9 (2.4–3.4)	3.1 (2.6–3.8)	0.022
CEA (ng/ml)	2.1 (1.1–4.6)	5.7 (2.7–30.7)	<0.001	2.4 (1.3–7.3)	5.7 (3.7–8.5)	<0.001
CA199 (U/ml)	7.0 (3.4–12.4)	30.3 (11.8–62.1)	<0.001	9.3 (5.5–20.1)	25.1 (8.7–119.5)	<0.001
<i>Pathological characteristics</i>						
Tumor size (cm), <5/≥5 [n (%)]	253/120 (67.8/32.2)	175/65 (72.9/27.1)	0.180	97/63 (60.6/39.4)	66/26 (71.7/28.3)	0.076
Location, rectal/colon [n (%)]	102/271 (27.3/72.7)	76/164 (31.7/68.3)	0.250	65/95 (40.6/59.4)	40/52 (43.5/56.5)	0.658
T stage, T1–T2/T3–T4 [n (%)]	93/280 (24.9/75.1)	78/162 (32.5/67.5)	0.041	70/90 (43.8/56.2)	19/73 (20.7/79.3)	<0.001
Lymph node metastasis, negative/positive [n (%)]	245/128 (65.7/34.3)	102/138 (42.5/57.5)	<0.001	91/69 (56.9/43.1)	27/65 (29.3/70.7)	<0.001
Tumor differentiation, poor/well [n (%)]	41/332 (11.0/89.0)	20/220 (8.3/91.7)	0.283	13/147 (8.1/91.9)	9/83 (9.8/90.2)	0.654

WBC white blood cell, Hb haemoglobin, RDW red cell distribution width, MCV mean corpuscular volume, MPV mean platelet volume, AGR albumin-to-globulin ratio, TBIL total bilirubin, DBIL direct bilirubin, TBA total bile acids, ALT alanine aminotransferase, AST aspartate aminotransferase, LDH lactate dehydrogenase, CK creatine kinase, Mb myoglobin, TG triglycerides, LDL-C low-density lipoprotein cholesterol, HDL-C high-density lipoprotein cholesterol, CEA carcinoembryonic antigen, CA199 carbohydrate antigen 199

CRLM-Lab6 prediction model construction

The discovery cohort was utilized for model construction. A total of five ML models were employed for performance comparison: Logistic regression, Linear SVC, Random forest, Decision tree, and SVM. Specifically, each model underwent training and parameter optimization using the training set, while the validation set was used to assess model performance. Among these algorithms, the Random forest

model demonstrated superior performance. The AUC of the Random forest model in the training set was 1.00 (Fig. 3a). Additionally, in the internal validation set, the AUC of the Random forest model was 0.93 (Fig. 3b), with a sensitivity of 0.82, specificity of 0.90, PPV of 0.83, and NPV of 0.88 (Table 2). Consequently, we selected the Random forest algorithm to construct the CRLM prediction model.

We further analyzed the relative importance of each variable in the Random forest model, ranking them in descending

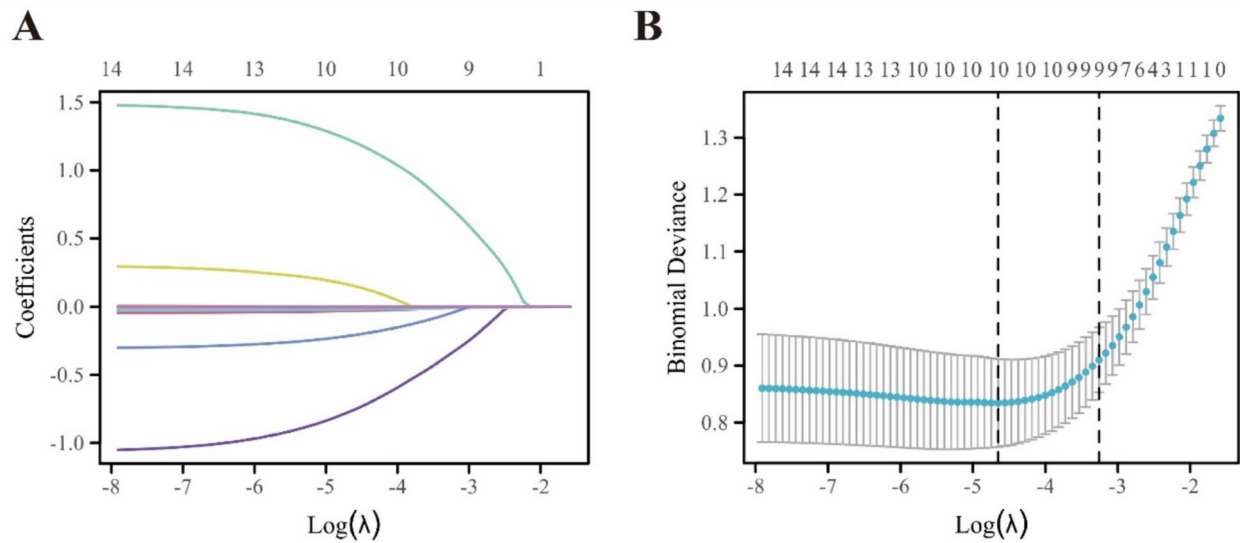


Fig. 2 Feature selection based on LASSO regression. **a** LASSO coefficient profiles of the fourteen risk factors. **b** Nine risk factors selected using LASSO regression analysis

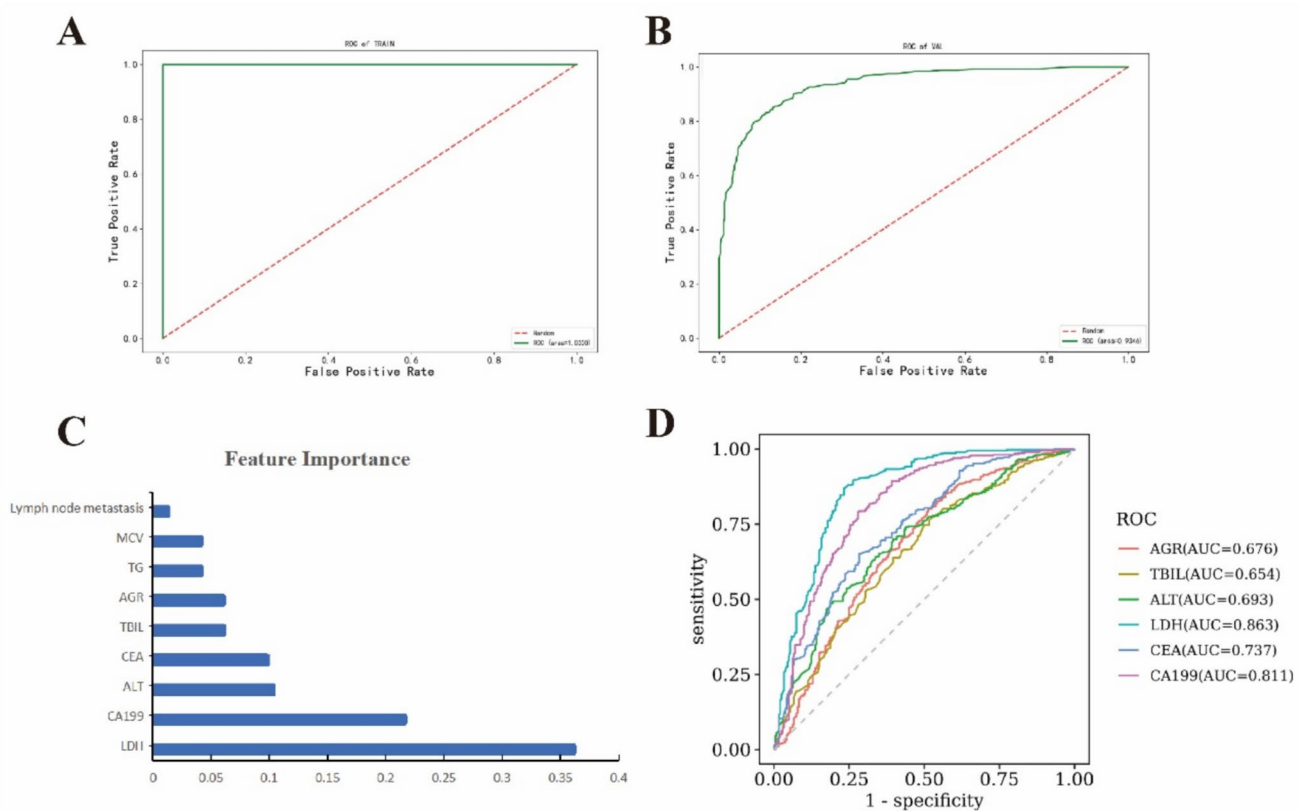


Fig. 3 ML-based prediction model for CRLM. **a** The ROC curve of the Random forest model in the internal training set. **b** The ROC curve of the Random forest model in the internal validation set. **c**

Weight of AGR, TBIL, ALT, LDH, CEA, CA199, TG, MCV and lymph node metastasis in the Random forest model. **d** ROC curve analysis of six features

Table 2 Performance of five models

Classifier	AUC	Sensitivity	Specificity	Accuracy	Precision	F1 score	PPV	NPV
Logistic regression	0.91	0.79	0.86	0.83	0.78	0.79	0.78	0.86
Linear SVC	0.91	0.78	0.87	0.84	0.80	0.79	0.80	0.86
Random forest	0.93	0.82	0.90	0.86	0.83	0.83	0.83	0.88
Decision tree	0.90	0.82	0.83	0.83	0.76	0.80	0.76	0.88
SVM	0.90	0.75	0.89	0.83	0.81	0.78	0.81	0.84

AUC area under the curve, *PPV* positive predictive value, *NPV* negative predictive value, *Linear SVC* Linear support vector classification, *SVM* support vector machine

order as follows: LDH, CA199, ALT, CEA, TBIL, AGR, TG, MCV, and lymph node metastasis (Fig. 3c). To simplify the model while retaining predictive power, we selected the top six variables by importance (LDH, CA199, ALT, CEA, TBIL, and AGR) to optimize the model. Consequently, a

refined Random forest model based on these six key variables was constructed and designated as CRLM-Lab6. Additionally, the ROC curve analysis demonstrated that these six features exhibited robust predictive performance for CRLM risk (Fig. 3d).

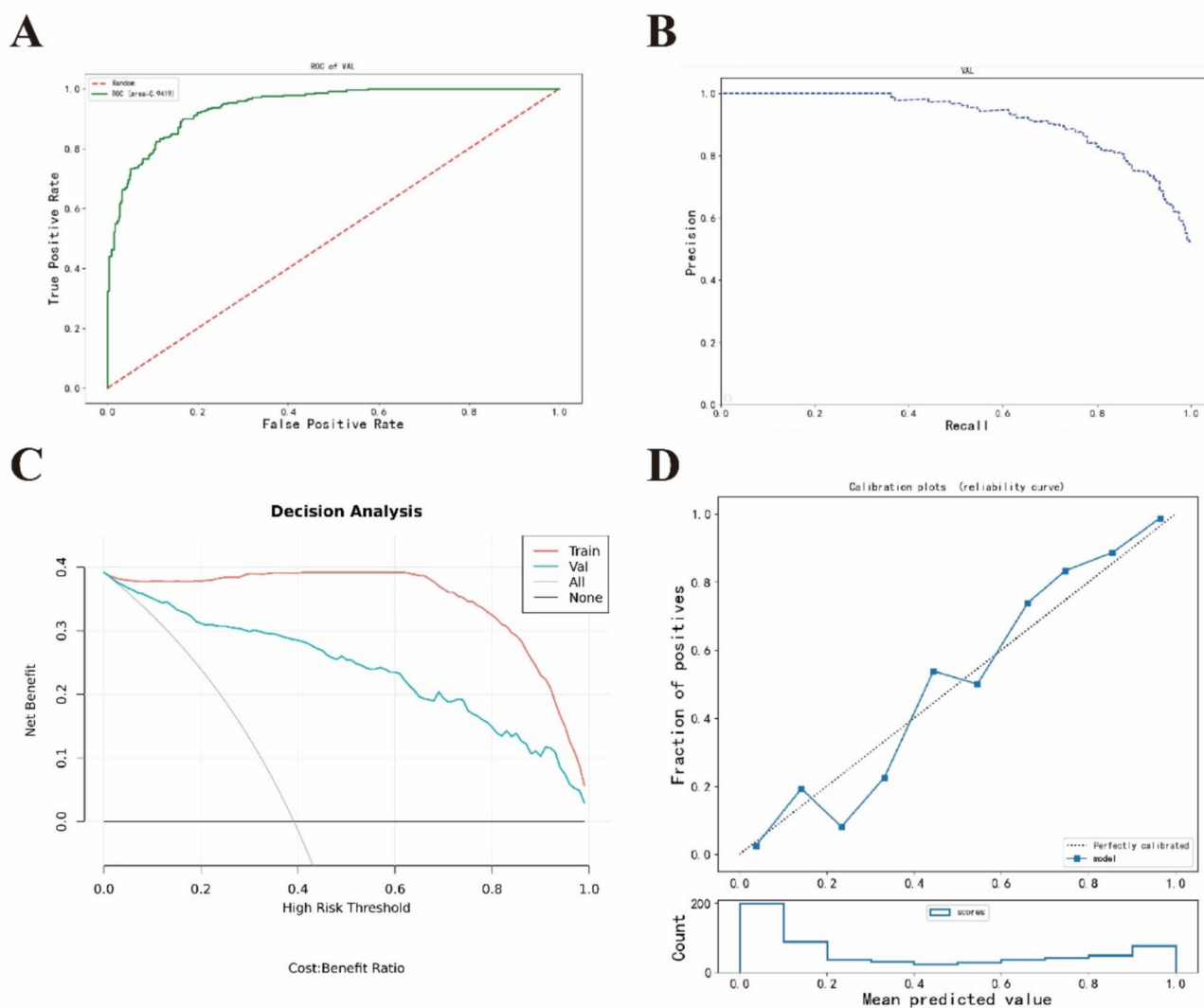


Fig. 4 Evaluation of the ability of the CRLM-Lab6 model to predict CRLM. **a** ROC curve analysis. **b** Precision recall curve analysis. **c** Decision curve analysis. **d** Calibration curve analysis

In the internal validation set, the CRLM-Lab6 model demonstrated exceptional performance with an AUC of 0.94, with a sensitivity of 0.88, specificity of 0.93 (Fig. 4a). The precision-recall curve indicated that the model achieved satisfactory performance (Fig. 4b). Additionally, the decision curve analysis confirmed that this model offers substantial clinical utility. (Fig. 4c). Furthermore, the calibration curve revealed a close alignment between the observed and predicted probabilities (Fig. 4d).

Validation of CRLM-Lab6 prediction model

To further evaluate the performance of the CRLM-Lab6 model, we conducted an analysis in a validation cohort comprising 252 subjects between January 2023 and September 2024, including 160 non-CRLM patients, and 92 CRLM patients. The ROC analysis demonstrated that the CRLM-Lab6 model effectively distinguished the CRLM patients from the non-CRLM patients, achieving an AUC of 0.96, with a sensitivity of 0.95 and specificity of 0.93.

A web page calculator of CRLM prediction model

Based on the CRLM-Lab6 model, we have developed an online calculator designed to predict the risk of CRLM occurrence. Users can input the relevant variables into the tool available at (<https://dxonline.deepwise.com/prediction/index.html?baseUrl=%2Fapi%2F&id=49967&topicName=undefined&from=share&platformType=wisdom>) to obtain a prediction of CRLM risk (Fig. 5).

Discussion

CRLM is a significant factor contributing to poor prognosis in CRC patients. Conventional imaging modalities, such as CT and MRI, lack the sensitivity for early detection of CRLM [21]. Consequently, there is an urgent need to develop more effective methods for identifying CRLM at an early stage. Early detection and accurate diagnosis of CRLM are critical for improving patient outcomes and



Fig. 5 The visualization of the CRLM-Lab6 model through the DxAI platform

survival rates. This study recruited a total of 865 participants, comprising 533 non-CRLM patients, and 332 CRLM patients. In the discovery cohort, we initially identified nine predictors using LASSO regression. We then evaluated the predictive performance of five ML algorithms: Logistic regression, Linear SVC, Random forest, Decision tree, and SVM. The Random forest model demonstrated superior performance with an AUC of 0.93. To simplify the model while maintaining its robustness, we developed the CRLM-Lab6 model based on the Random forest algorithm. In the external validation set, the CRLM-Lab6 model exhibited stable and excellent discriminatory power in the non-CRLM, achieving AUC of 0.96. Furthermore, we created an online calculator to predict the risk of CRLM occurrence. By inputting relevant variables, clinicians can estimate the probability of CRLM. This tool is particularly valuable in regions with limited medical resources, aiding clinicians in early identification of CRLM.

In this study, we developed a CRLM prediction model using blood routine and biochemical test data. This model aids clinicians in early diagnosis of CRLM and facilitates the formulation of personalized treatment strategies. The model incorporates six common clinical indicators: LDH, CA199, ALT, CEA, TBIL, and AGR. Research has demonstrated that CRLM can lead to abnormal liver function, often associated with hepatocyte damage, resulting in elevated levels of enzymes such as ALT and AST, as well as bilirubin [22–25]. LDH, an enzyme critical for cellular metabolism, is increased in CRLM patients due to heightened tumor metabolic activity and the resultant metabolic burden on the liver, leading to its release from hepatocytes into the bloodstream [26, 27]. AGR, an indicator of systemic inflammation, is strongly correlated with tumor severity and poor prognosis [28, 29]. Tumor markers such as CEA and CA199 are crucial for preoperative screening of CRC and predicting post-surgical recurrence and metastasis [30]. For instance, studies have shown that CEA and CA199 expression levels correlate with factors like tumor invasion depth, diameter, and metastasis in CRC [31].

Furthermore, our model outperforms those presented in the existing literature. For instance, wang et al. [32]. constructed a CRLM model combining the albumin-bilirubin score, ALT levels, and CEA, achieving an AUC of 0.92, sensitivity of 0.78, and specificity of 0.95. A detection method based on circulating tumor DNA (ctDNA) achieves an AUC value as high as 0.90 in predicting CRLM [33]. Nevertheless, its broad application remains constrained to some extent by factors such as tumor heterogeneity, the dynamic nature of ctDNA, technical complexity, and relatively high detection costs. Overall, the multi-biomarker prediction model established in this study utilizes readily available hematological data, offering significant cost-effectiveness and aiding in the early diagnosis and treatment of CRLM.

This study has several limitations. First, the cross-sectional design of this study restricts our ability to determine the temporal relationship between routine laboratory test data and tumor metastasis. To overcome this limitation, future studies will involve longitudinal follow-up of CRC patients to validate the model's accuracy. Second, as a single-center retrospective study, potential biases in case selection may exist. Specifically, resource constraints limit this study to a single hospital setting. In subsequent research, we plan to incorporate data from additional centers to enhance the generalizability of the model. Third, while the model demonstrates high accuracy in differentiating CRLM, technical and resource limitations precluded us from integrating it with other clinical practices, such as imaging data, for further refinement. This aspect will be explored in future analyses.

Conclusion

This study established a CRLM prediction model using ML algorithms. The Random forest algorithm was employed to construct the model named CRLM-Lab6, which incorporates six variables: LDH, CA199, ALT, CEA, TBIL, and AGR. This model demonstrated high predictive accuracy and clinical utility, thereby assisting clinicians in early identification of CRLM patients.

Author contribution SF Writing—original draft, Methodology, Investigation, Formal analysis, Conceptualization. MZ: Writing—review and editing, Writing—original draft, Conceptualization. ZH Supervision, Methodology. XX Supervision, Methodology. BZ Writing—review and editing, Supervision, Funding acquisition, Conceptualization. SSF. and MLZ. wrote the main manuscript text, ZXH. and XMX. prepared Figs. 1–3. All authors reviewed the manuscript.

Funding This study was supported by Hunan Province Key Research and Development Project (2024JK2107), and the Natural Science Foundation of Hunan Province (2022JJ30987).

Data availability No datasets were generated or analysed during the current study.

Declarations

Conflicts of interest The authors declare no competing interests.

Ethical approval Ethical approval was waived by the local Ethics Committee of Xiangya Hospital, Central South University in view of the retrospective nature of the study and all the procedures being performed were part of the routine care.

Consent for publication Not applicable.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit

to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

References

- Bray F, Laversanne M, Sung H, et al. Global cancer statistics 2022: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *Ca Cancer J Clin*. 2024;74:229–63.
- Cheng XF, Zhao F, Chen D, et al. Current landscape of preoperative neoadjuvant therapies for initial resectable colorectal cancer liver metastasis. *World J Gastroenterol*. 2024;30:663–72.
- Burnett-Hartman AN, Lee JK, Demb J, et al. An update on the epidemiology, molecular characterization, diagnosis, and screening strategies for early-onset colorectal cancer. *Gastroenterology*. 2021;160:1041–9.
- Filoni E, Musci V, Di Rito A, et al. Multimodal management of colorectal liver metastases: state of the art. *Oncol Rev*. 2023;17:11799.
- Itenberg ER, Lozano AM. Surgical and interventional management of liver metastasis. *Clin Colon Rectal Surg*. 2024;37:80–4.
- Wang Y, Zhong X, He X, et al. Liver metastasis from colorectal cancer: pathogenetic development, immune landscape of the tumour microenvironment and therapeutic approaches. *J Exp Clin Cancer Res*. 2023;42:177.
- Liu QL, Zhou H, Zhou ZG, et al. Colorectal cancer liver metastasis: genomic evolution and crosstalk with the liver microenvironment. *Cancer Metastasis Rev*. 2023;42:575–87.
- Datta J, Narayan RR, Kemeny NE, et al. Role of hepatic artery infusion chemotherapy in treatment of initially unresectable colorectal liver metastases: a review. *Jama Surg*. 2019;154:768–76.
- Xu L, Cai S, Cai G, et al. Imaging diagnosis of colorectal liver metastases. *World J Gastroenterol*. 2011;17:4654–9.
- Jones A, Findlay A, Knight SR, et al. Follow up after surgery for colorectal liver metastases: a systematic review. *Eur J Surg Oncol*. 2023;49: 107103.
- Olaker VR, Fry S, Terebuh P, et al. With big data comes big responsibility: strategies for utilizing aggregated, standardized, de-identified electronic health record data for research. *Clin Transl Sci*. 2025;18: e70093.
- Coronado GD, Nielson CM, Keast EM, Petrik AF, Suls JM. The influence of multi-morbidities on colorectal cancer screening recommendations and completion. *Cancer Causes Control*. 2021;32:555–65.
- Parwaiz I, Hakeem A, Nwogwugwu O, et al. Does ALT Correlate with survival after liver resection for colorectal liver metastases? *J Clin Exp Hepatol*. 2022;12:1285–92.
- Joechle K, Goumard C, Vega EA, et al. Long-term survival after post-hepatectomy liver failure for colorectal liver metastases. *Hpb (Oxford)*. 2019;21:361–9.
- Waite S, Davenport MS, Graber ML, et al. Opportunity and opportunism in artificial intelligence-powered data extraction: a value-centered approach. *Ajr Am J Roentgenol*. 2024;204:1–11.
- Swinkels L, Bennis FC, Ziesemer KA, et al. The use of deep learning and machine learning on longitudinal electronic health records for the early detection and prevention of diseases: scoping review. *J Med Internet Res*. 2024;26: e48320.
- Greener JG, Kandathil SM, Moffat L, et al. A guide to machine learning for biologists. *Nat Rev Mol Cell Biol*. 2022;23:40–55.
- Jiang T, Gradus JL, Rosellini AJ. Supervised machine learning: a brief primer. *Behav Ther*. 2020;51:675–87.
- Miller-Atkins G, Acevedo-Moreno LA, Grove D, et al. Breath metabolomics provides an accurate and noninvasive approach for screening cirrhosis, primary, and secondary liver tumors. *Hepatol Commun*. 2020;4:1041–55.
- Amin MB, Greene FL, Edge SB, et al. The eighth edition AJCC cancer staging manual: continuing to build a bridge from a population-based to a more personalized approach to cancer staging. *Ca Cancer J Clin*. 2017;67:93–9.
- Zhou H, Liu Z, Wang Y, et al. Colorectal liver metastasis: molecular mechanism and interventional therapy. *Signal Transduct Target Ther*. 2022;7:70.
- Zou Y, Zhong L, Hu C, et al. Association between the alanine aminotransferase/aspartate aminotransferase ratio and new-onset non-alcoholic fatty liver disease in a nonobese Chinese population: a population-based longitudinal study. *Lipids Health Dis*. 2020;19:245.
- Lam C, Bharwani AA, Chan E, et al. A machine learning model for colorectal liver metastasis post-hepatectomy prognostications. *Hepatobiliary Surg Nutr*. 2023;12:495–506.
- Liu X, Meng J, Xu H, et al. Alpha-fetoprotein to transaminase ratio is related to higher diagnostic efficacy for hepatocellular carcinoma. *Medicine (Baltimore)*. 2019;98: e15414.
- Scheipner L, Smolle MA, Barth D, et al. The AST/ALT ratio is an independent prognostic marker for disease-free survival in stage II and III colorectal carcinoma. *Anticancer Res*. 2021;41:429–36.
- Bai L, Lin ZY, Lu YX, et al. The prognostic value of preoperative serum lactate dehydrogenase levels in patients underwent curative-intent hepatectomy for colorectal liver metastases: a two-center cohort study. *Cancer Med*. 2021;10:8005–19.
- Bai L, Yan XL, Lu YX, et al. Circulating lipid- and inflammation-based risk (CLIR) score: a promising new model for predicting outcomes in complete colorectal liver metastases resection. *Ann Surg Oncol*. 2022;29:4308–23.
- Li J, Zhu N, Wang C, et al. Preoperative albumin-to-globulin ratio and prognostic nutritional index predict the prognosis of colorectal cancer: a retrospective study. *Sci Rep*. 2023;13:17272.
- Lv GY, An L, Sun XD, Hu YL, Sun DW. Pretreatment albumin to globulin ratio can serve as a prognostic marker in human cancers: a meta-analysis. *Clin Chim Acta*. 2018;476:81–91.
- Paredes AZ, Hyer JM, Tsimiligras DI, et al. A novel machine-learning approach to predict recurrence after resection of colorectal liver metastases. *Ann Surg Oncol*. 2020;27:5139–47.
- Polat E, Duman U, Duman M, et al. Diagnostic value of preoperative serum carcinoembryonic antigen and carbohydrate antigen 19–9 in colorectal cancer. *Curr Oncol*. 2014;21:e1–7.
- Wang ZM, Pan SP, Zhang JJ, et al. Prediction and analysis of albumin-bilirubin score combined with liver function index and carcinoembryonic antigen on liver metastasis of colorectal cancer. *World J Gastrointest Surg*. 2024;16:1670–80.
- Wang XY, Zhang R, Han JH, et al. Early circulating tumor DNA dynamics predict neoadjuvant therapy response and recurrence in colorectal liver metastases: a prospective study. *Ann Surg Oncol*. 2023;30(8):5252–63.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.