

# Survival Prediction Model for Patients with Hepatocellular Carcinoma and Extrahepatic Metastasis Based on XGBoost Algorithm

Jihye Lim<sup>1</sup>, Hyeon-Gi Jeon<sup>2</sup>, Yeonjoo Seo<sup>1</sup>, Moonjin Kim<sup>3</sup>, Ja Un Moon<sup>4</sup>, Se Hyun Cho<sup>1</sup>

<sup>1</sup>Division of Gastroenterology and Hepatology, Department of Internal Medicine, Yeouido St. Mary's Hospital, College of Medicine, The Catholic University of Korea, Seoul, Republic of Korea; <sup>2</sup>Department of Core Platform Team, SOCAR Incorporated, Seoul, Republic of Korea; <sup>3</sup>Department of Internal Medicine, Yeouido St. Mary's Hospital, College of Medicine, The Catholic University of Korea, Seoul, Korea; <sup>4</sup>Department of Pediatrics, Yeouido St. Mary's Hospital, College of Medicine, The Catholic University of Korea, Seoul, Republic of Korea

Correspondence: Se Hyun Cho, Division of Gastroenterology and Hepatology, Department of Internal Medicine, Yeouido St. Mary's Hospital, College of Medicine, The Catholic University of Korea, 63-ro 10, Yeongdeungpo-gu, Seoul, 07345, Republic of Korea, Tel +82-2-3779-1269, Fax +82-02-780-3132, Email chowhang@catholic.ac.kr

**Purpose:** Accurate estimation of survival is of utmost importance in patients with hepatocellular carcinoma (HCC) and extrahepatic metastasis. This study aimed to develop a survival prediction model using real-world data.

**Patients and Methods:** A total of 993 patients with treatment-naïve HCC and extrahepatic metastasis were included from 13 Korean hospitals between 2013 and 2018. Patients were randomly divided into a training set (70.0%) and a test set (30.0%). The eXtreme Gradient Boosting (XGBoost) algorithm was applied to predict survival at 3, 6, and 12 months.

**Results:** The mean age of the patients was  $60.8 \pm 12.3$  years, and 85.4% were male. During the study period, 96.1% died, and median survival duration was 4.0 months. In multivariate analysis, Child-Pugh class, number and size of tumors, presence of vascular or bile duct invasion, lung or bone metastasis, serum AFP, and primary anti-HCC treatment were associated with survival. We constructed a model for survival prediction based on the relevant variables, which is available online (<https://metastatic-hcc.onrender.com/form>). Our model demonstrated high performance, with areas under the receiver operating characteristic curves of 0.778, 0.794, and 0.784 at 3, 6, and 12 months, respectively. Feature importance analysis indicated that the primary anti-HCC treatment had the highest importance.

**Conclusion:** We developed a model to predict the survival of patients with HCC and extrahepatic metastasis, which demonstrated good discriminative ability. Our model would be helpful for personalized treatment and for improving the prognosis.

**Keywords:** liver neoplasms, prognosis, survival rate, probability, algorithms

## Introduction

Hepatocellular carcinoma (HCC) is a multifaceted disease that requires careful consideration of treatment decisions. Unlike other solid cancers which are mostly treated according to the tumor burden, HCC treatment requires careful evaluation of concomitant liver disease, liver function, and general performance status as well as intra- and extrahepatic tumor burden.<sup>1-4</sup> Also, HCCs have inter-patient and intra-tumoral heterogeneity arising from chronic liver inflammation with complex pathogenesis, accompanied by a range of genetic and epigenetic changes.<sup>5</sup> Thereby, patients with HCC and extrahepatic metastasis, classified under the same stage, may exhibit varied clinical courses with estimated survival duration of 7–16 months.<sup>2,4,6</sup> However, survival prediction of patients with HCC and extrahepatic metastasis often relies on a few studies analyzed with a limited set of factors, or on physicians' experience, which is prone to inaccurate. Hence, there has been an urgent need for systematic and personalized survival prediction for these patients.

The eXtreme Gradient Boosting (XGBoost) algorithm, a product of modern information technology, was introduced by Chen and Guestrin in 2016.<sup>7</sup> This algorithm employs ensemble learning, one of the most widely used machine learning methods, to construct a single generalized model with high predictive ability. It leverages the gradient descent method to generate decision trees for classification and regression, thereby enhancing robustness.<sup>8</sup> XGBoost is renowned

for improving the accuracy and computational speed of machine learning algorithms, and has found widespread applications in various fields, including the medical domain.<sup>7–10</sup> Another strength of XGBoost is its interpretability. The decision-making process of XGBoost is understandable, unlike other machine learning algorithms.<sup>11,12</sup> Utilizing XGBoost in the context of patients with HCC and extrahepatic metastasis can enhance survival prediction performance.

This study aimed to develop a model for predicting the survival of patients with HCC and extrahepatic metastasis using the XGBoost algorithm with risk factors associated with survival.

## Materials and Methods

### Study Design and Population

Between January 2013 and December 2018, 9083 patients initially diagnosed with HCC and extrahepatic metastasis were recruited from the National Cancer Registry in Korea. The National Cancer Center selected 13 hospitals, of which over 75% of patients with liver cancer in Korea were treated to build the registry. A systematic extraction method was used to sample 10% of the initially diagnosed HCC patients, and their comprehensive medical data were investigated by medical record administrators. Death status and date of death were obtained from the Korean Statistical Office. The diagnosis of HCC was based on radiological hallmarks in multiphase CT or MRI, arterial phase enhancement, and portal or delayed phase washout appearance, according to the KLCA-NCC Korea practice guidelines.<sup>2</sup> Of these patients, we excluded the following: 1) 255 patients with missing data; 2) 243 patients who died within 30 days after initial HCC diagnosis; 3) 8585 patients without extrahepatic metastasis. The 993 patients were finally enrolled (Figure 1). This study was approved by the Institutional Review Board (IRB) of Catholic Medical Center (IRB No. SC22ZISE0092), and the need for informed consent was waived owing to the use of de-identified data. It was conducted in accordance with the Declaration of Helsinki and Istanbul.

### Data Collection

We gathered data on various clinical variables at the time of initial HCC diagnosis that are known to impact prognosis: 1) baseline patient-related variables including age, sex, body mass index (BMI), alcohol consumption, smoking, presence of diabetes, hypertension, dyslipidemia, and Eastern Cooperative Oncology Group (ECOG) performance status; 2) liver-related variables obtained for the etiology of liver disease, Child-Pugh class, Model for End-stage Liver Disease (MELD) score, and cirrhosis by image or fibrosis-4 (FIB-4) index  $\geq 2.67$ ;<sup>13</sup> 3) laboratory findings including platelet, prothrombin time, aspartic acid transaminase (AST), alanine transaminase (ALT), albumin, total bilirubin, creatinine, sodium, glucose, and total cholesterol; 4) tumor-related factors including number of tumors, size of the tumor, presence of intrahepatic vascular or bile duct invasion, extrahepatic metastatic site classified as regional lymph node, lung, bone, distal lymph node and others, alpha-fetoprotein (AFP), and initial anti-HCC treatment.

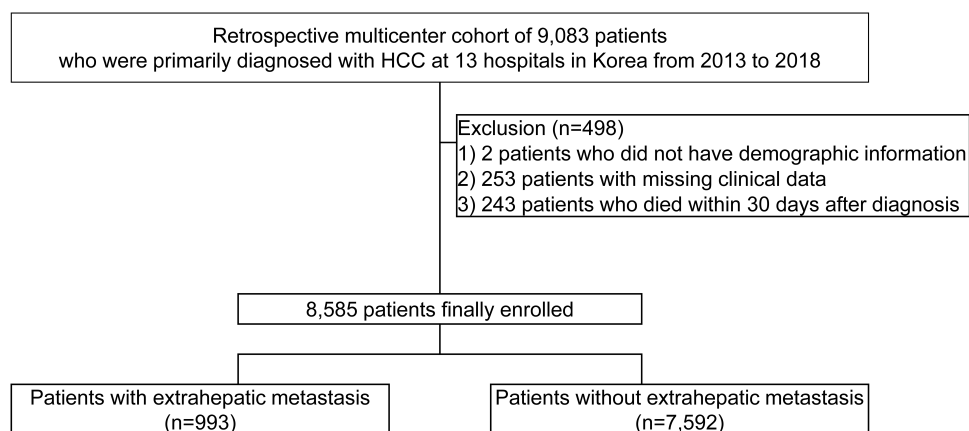


Figure 1 Study flow.

## Outcomes

The primary outcome of this study was to identify the clinical factors that affect survival and to estimate the survival of patients with HCC and extrahepatic metastasis at 3, 6, and 12 months after the initial diagnosis. In addition, we sought to develop a personalized model that could accurately predict survival.

## Statistical Analysis and Construction of a Prediction Model

Categorical variables were expressed as numbers with percentages and analyzed using the chi-square or Fisher's exact test. Continuous variables were summarized as mean  $\pm$  standard deviation (SD) and analyzed using the *t*-test or Mann–Whitney *U*-test. Survival along with clinical variables were estimated and compared with the Kaplan–Meier and Log rank test. Traditional Cox proportional hazard models were used to identify factors that could significantly influence survival. The hazard ratios (HR) and 95% confidence intervals (CI) were summarized. Statistical significance was set at  $P < 0.05$ .

In addition, we used XGBoost, a gradient-boosted decision tree algorithm, to build a prediction model with complete case analysis. We utilized multiple imputations to estimate missing data with less than 10% of patients who did not introduce bias.<sup>14</sup> Variables with more than 10% of missing data were excluded for model development.

The data was randomly split into a training set and a test set in a 7:3 ratio. The training set was used for model training, and the test set was employed for validation to assess the model's accuracy. We selected variables that demonstrated statistical significance in multivariable analysis or clinical relevance to include in the model. The response variables, representing survival information at each time point, were obtained.

To build our prediction models, we used the following parameters: number of estimators = 300, maximum depth of each tree = 15, *colsample.bytree* = 0.75, and *gamma* = 3.87. We utilized ensemble learning methods such as Bootstrap Aggregating and Boosting to integrate decision trees, effectively reducing both bias-related and variance-related errors.<sup>15</sup> The results generated by XGBoost ranged between 0 and 1, representing the probability of survival. Patients with results greater than 0.5 were classified as dead, while those with results less than 0.5 were classified as alive.

The performance of model was assessed by evaluating its accuracy, precision (positive predictive value), recall (sensitivity), F1 score (harmonic mean of precision and recall), and area under the receiver operating characteristic curve (AUC). An AUC value greater than 0.5 is considered statistically significant in evaluating binary classifiers.<sup>16</sup> Additionally, we analyzed the feature importance and decision trees of each model to elucidate and understand the decision-making process. Finally, we utilized a web-hosting service named “Render” to upload our model and make it available on a website for broader accessibility.<sup>17</sup> All statistical analyses were performed with R software version 4.1.3, SAS version 9.4 (SAS Institute, Inc., Cary, NC, USA), and Python 3.9.16 (Python Software Foundation, Delaware, USA).

## Results

### Characteristics of the Study Population

Of the 993 patients, 695 (70.0%) and 298 (30.0%) were assigned to the training and test sets, respectively. The baseline characteristics of the two groups are shown in Table 1. The mean ( $\pm$  SD) age was 60.8 ( $\pm$  12.3) years. Males accounted for 85.4% and hepatitis B virus (HBV) infection was the most common cause of HCC (60.1%). Among the study participants, 56.1% had Child-Pugh class A liver function, 48.7% had more than five masses, and 48.9% had the tumors larger than 10 cm in diameter. Intrahepatic vascular or bile duct invasion was observed in 63.3% of the patients. Regional lymph nodes were the most frequent metastatic sites (51.5%), followed by the lungs (39.8%), distant lymph nodes (21.1%), bones (20.0%), and others (17.9%). The mean AFP was 44,108.7  $\pm$  212,265.8 ng/dL. Supportive care (36.8%) was the most common first-line treatment, followed by systemic chemotherapy (28.8%), and transarterial therapy (24.8%). There was no statistically significant difference in the baseline characteristics between the two datasets.

### Survival Analysis and Risk Factors Associated with Mortality

During a median follow-up of 4.0 months (range, 2.00–9.02 months), 96.1% of the patients died in our study cohort. Survival rates were 62.5%, 38.6%, 26.8%, 21.7%, and 17.8% at 3, 6, 9, 12, and 15 months, respectively.

**Table I** Baseline Characteristics of the Study Population

	<b>Total (N=993)</b>	<b>Train set (N=695)</b>	<b>Test set (N=298)</b>	<b>P-value</b>
<i>Patients-related factors</i>				
Age, years	60.8 ± 12.3	60.4 ± 12.5	61.6 ± 11.6	0.185
Male sex, n (%)	848 (85.4%)	587 (84.5%)	261 (87.6%)	0.238
Alcohol	418 (42.1%)	281 (40.4%)	137 (46.0%)	0.121
Smoking	516 (52.0%)	362 (52.1%)	154 (51.7%)	0.961
Diabetes	262 (26.4%)	177 (25.5%)	85 (28.5%)	0.356
Hypertension	324 (32.6%)	226 (32.5%)	98 (32.9%)	0.969
Dyslipidemia	722 (72.7%)	506 (72.8%)	216 (72.5%)	0.207
Body mass index, kg/m <sup>2</sup>	23.2 ± 3.5	23.2 ± 3.5	23.0 ± 3.7	0.407
ECOG performance status				0.918
0	337 (33.9%)	241 (34.7%)	96 (32.2%)	
1	240 (24.2%)	165 (23.7%)	75 (25.2%)	
2	60 (6.0%)	41 (5.9%)	19 (6.4%)	
3	27 (2.7%)	19 (2.7%)	8 (2.7%)	
4	15 (1.5%)	12 (1.7%)	3 (1.0%)	
Missing	314 (31.6%)	217 (31.2%)	97 (32.6%)	
<i>Liver-related factors</i>				
Etiology of liver disease				0.080
HBV infection	597 (60.1%)	424 (61.0%)	173 (58.1%)	
HCV infection	77 (7.8%)	48 (6.9%)	29 (9.7%)	
Alcohol	148 (14.9%)	95 (13.7%)	53 (17.8%)	
Others	171 (17.2%)	128 (18.4%)	43 (14.4%)	
Child-Pugh class				0.161
A	557 (56.1%)	382 (55.0%)	175 (58.7%)	
B	367 (37.0%)	258 (37.1%)	109 (36.6%)	
C	69 (6.9%)	55 (7.9%)	14 (4.7%)	
MELD score	8.2 ± 1.8	8.3 ± 1.9	8.1 ± 1.6	0.152
Cirrhosis				0.313
Without cirrhosis	146 (14.7%)	101 (14.5%)	45 (15.1%)	
With cirrhosis	300 (30.2%)	220 (31.7%)	80 (26.8%)	
Missing	547 (55.1%)	374 (53.8%)	173 (58.1%)	
<i>Laboratory findings</i>				
Platelets, 1000/mm <sup>3</sup>	206.0 ± 110.5	208.8 ± 118.1	199.4 ± 90.0	0.171
Prothrombin time, INR	1.2 ± 0.3	1.2 ± 0.3	1.2 ± 0.2	0.113
AST, IU/L	123.6 ± 118.5	126.3 ± 122.6	116.9 ± 107.7	0.345
ALT, IU/L	61.6 ± 57.0	62.5 ± 55.6	59.6 ± 60.4	0.475
Albumin, g/dL	3.5 ± 0.6	3.5 ± 0.6	3.6 ± 0.7	0.234
Total bilirubin, mg/dL	2.1 ± 3.5	2.1 ± 3.6	1.9 ± 3.4	0.280
Creatinine, mg/dL	0.9 ± 0.6	1.0 ± 0.7	0.9 ± 0.4	0.375
Sodium, mmol/L	136.4 ± 4.4	136.4 ± 4.2	136.4 ± 4.7	0.897
Glucose, mg/dL	131.1 ± 62.2	130.6 ± 61.5	132.4 ± 63.7	0.724
Total cholesterol, mg/dL	171.2 ± 66.7	171.4 ± 66.1	170.9 ± 68.0	0.924
<i>Tumor-related factors</i>				
Number of tumors				0.102
1	386 (38.9%)	258 (37.1%)	128 (43.0%)	
2	87 (8.8%)	55 (7.9%)	32 (10.7%)	
3	26 (2.6%)	17 (2.4%)	9 (3.0%)	

(Continued)

Table 1 (Continued).

	Total (N=993)	Train set (N=695)	Test set (N=298)	P-value
4	10 (1.0%)	8 (1.2%)	2 (0.7%)	
≥5	484 (48.7%)	357 (51.4%)	127 (42.6%)	
Size of tumor, cm				0.472
≤ 2	48 (4.8%)	37 (5.3%)	11 (3.7%)	
2–5	134 (13.5%)	89 (12.8%)	45 (15.1%)	
5–7	117 (11.8%)	78 (11.2%)	39 (13.1%)	
7–10	208 (20.9%)	152 (21.9%)	56 (18.8%)	
> 10	486 (48.9%)	339 (48.8%)	147 (49.3%)	
Vascular or bile duct invasion	629 (63.3%)	444 (63.9%)	185 (62.1%)	0.639
Extrahepatic involvement				
Regional lymph node	511 (51.5%)	354 (50.9%)	157 (52.7%)	0.663
Lung	395 (39.8%)	280 (40.3%)	115 (38.6%)	0.667
Bone	199 (20.0%)	138 (19.9%)	61 (20.5%)	0.893
Distant lymph node	210 (21.1%)	147 (21.2%)	63 (21.1%)	>0.999
Others	178 (17.9%)	124 (17.8%)	54 (18.1%)	0.988
AFP, ng/dL	44108.7 ± 212,265.8	41,953.1 ± 204,184.4	49,136.1 ± 230,296.4	0.642
Primary anti-HCC treatment				0.072
Surgical resection	42 (4.2%)	24 (3.5%)	18 (6.0%)	
Liver transplantation	3 (0.3%)	2 (0.3%)	1 (0.3%)	
Local ablation therapy	7 (0.7%)	3 (0.4%)	4 (1.3%)	
Transarterial therapy	246 (24.8%)	163 (23.5%)	83 (27.9%)	
Chemotherapy	286 (28.8%)	202 (29.1%)	84 (28.2%)	
Radiation therapy	44 (4.4%)	36 (5.2%)	8 (2.7%)	
Supportive care	365 (36.8%)	265 (38.1%)	100 (33.6%)	

**Notes:** Values are expressed as the mean ± standard deviation or frequency (percentage).

**Abbreviations:** AFP, alpha-fetoprotein; ALT, alanine aminotransferase; AST, aspartate aminotransferase; ECOG, Eastern Cooperative Oncology Group; HBV, hepatitis B virus; HCV, hepatitis C virus; INR, international normalized ratio; MELD, model for end-stage liver disease.

In a multivariate analysis, Child-Pugh class (class B: HR, 1.52 [95% CI, 1.31–1.77],  $P < 0.001$ ; class C: HR, 1.60 [95% CI, 1.18–2.17],  $P = 0.003$ ), platelets,  $\log_{10}$  1000/mm<sup>3</sup> (HR, 1.55 [95% CI, 1.14–2.11],  $P = 0.005$ ), sodium, mmol/L (HR, 0.97 [95% CI, 0.95–0.99],  $P < 0.001$ ), number of tumors ( $\geq 5$ : HR, 1.35 [95% CI, 1.17–1.56],  $P < 0.001$ ), size of tumors (5–7 cm: HR, 1.48 [95% CI, 1.01–2.17],  $P = 0.046$ ; 7–10 cm: HR, 1.64 [95% CI, 1.14–2.35],  $P = 0.008$ ;  $> 10$  cm: HR, 1.71 [95% CI, 1.2–2.45],  $P = 0.003$ ), presence of vascular or bile duct invasion (HR, 1.17 [95% CI, 1.01–1.36],  $P = 0.035$ ), lung metastasis (HR, 1.19 [95% CI, 1.01–1.39],  $P = 0.032$ ), bone metastasis (HR, 1.36 [95% CI, 1.14–1.62],  $P = 0.001$ ), serum AFP,  $\log_{10}$  ng/dL (HR, 1.08 [95% CI, 1.03–1.14],  $P = 0.001$ ), and primary anti-HCC treatment other than supportive care (surgical resection: HR 0.24 [95% CI, 0.16–0.35],  $P < 0.001$ ; liver transplantation: HR, 0.28 [95% CI, 0.09–0.89],  $P = 0.031$ ; local ablation therapy: HR, 0.35 [95% CI, 0.14–0.87],  $P = 0.024$ ; transarterial therapy: HR, 0.46 [95% CI, 0.39–0.56],  $P < 0.001$ ; chemotherapy: HR, 0.67 [95% CI, 0.57–0.8],  $P < 0.001$ ; radiation therapy: HR, 0.61 [95% CI, 0.44–0.84],  $P = 0.003$ ) were associated with mortality (Table 2 and [Supplementary Figure 1](#)).

## Model Construction and Evaluation

Based on survival analysis, the following variables were included to build a model: age, sex, BMI, etiology of liver disease, Child-Pugh class, platelet count, ALT level, sodium level, number of tumors, tumor size, presence of vascular or bile duct invasion, metastatic site, serum AFP level, and primary treatment method for HCCs.

We built survival prediction model for 3, 6, and 12-months by using the training dataset. The performances of these models are summarized in Table 3. Figure 2 shows AUC of each model. The AUC values of the test set at 3, 6, and 12-months were

**Table 2** Univariate and Multivariate Analyses for Overall Survival

	Univariate Analysis HR (95% CI)	P-value	Multivariate Analysis HR (95% CI)	P-value
<i>Patients-related factors</i>				
Age, years	1.00 (1.00–1.01)	0.335		
Male sex, n (%)	1.11 (0.93–1.33)	0.265		
Alcohol	0.88 (0.77–1.00)	0.054		
Smoking	1.01 (0.89–1.15)	0.858		
Diabetes	0.98 (0.85–1.13)	0.742		
Hypertension	1.00 (0.87–1.15)	0.995		
Dyslipidemia	1.00 (1.00–1.00)	0.128		
Body mass index, kg/m <sup>2</sup>	0.98 (0.96–1.00)	0.043	1.00 (0.98–1.02)	0.935
<i>Liver-related factors</i>				
Etiology of liver disease				
HBV infection	I (reference)			
HCV infection	1.11 (0.87–1.41)	0.411		
Alcohol	0.83 (0.69–1.00)	0.054		
Others	1.06 (0.89–1.26)	0.49		
Child-Pugh class				
A	I (reference)			
B	1.72 (1.50–1.97)	< 0.001	1.52 (1.31–1.77)	< 0.001
C	2.64 (2.04–3.40)	< 0.001	1.60 (1.18–2.17)	0.003
<i>Laboratory findings</i>				
Platelets, log <sub>10</sub> 1000/mm <sup>3</sup>	1.58 (1.18–2.11)	0.002	1.55 (1.14–2.11)	0.005
ALT, log <sub>10</sub> IU/L	1.39 (1.15–1.68)	0.001	0.93 (0.76–1.13)	0.450
Creatinine, mg/dL	1.06 (0.95–1.17)	0.297		
Sodium, mmol/L	0.94 (0.93–0.95)	< 0.001	0.97 (0.95–0.99)	< 0.001
<i>Tumor-related factors</i>				
Number of tumors				
I	I (reference)		I (reference)	
2	1.10 (0.86–1.39)	0.455	1.12 (0.88–1.44)	0.356
3	0.75 (0.50–1.14)	0.183	1.14 (0.74–1.75)	0.545
4	0.80 (0.43–1.51)	0.494	0.78 (0.41–1.48)	0.445
≥5	1.62 (1.41–1.86)	< 0.001	1.35 (1.17–1.56)	< 0.001
Size of tumor, cm				
≤ 2	I (reference)		I (reference)	
2–5	1.55 (1.09–2.21)	0.015	1.25 (0.87–1.81)	0.226
5–7	1.89 (1.32–2.71)	0.001	1.48 (1.01–2.17)	0.046
7–10	2.23 (1.59–3.13)	< 0.001	1.64 (1.14–2.35)	0.008
> 10	2.89 (2.10–3.99)	< 0.001	1.71 (1.20–2.45)	0.003
Vascular or bile duct invasion	1.49 (1.31–1.70)	< 0.001	1.17 (1.01–1.36)	0.035
Extrahepatic involvement				
Regional lymph node	0.92 (0.81–1.05)	0.215	1.12 (0.97–1.30)	0.130
Lung	1.53 (1.34–1.74)	< 0.001	1.19 (1.01–1.39)	0.032
Bone	1.14 (0.98–1.34)	0.096	1.36 (1.14–1.62)	0.001
Distant lymph node	1.22 (1.05–1.43)	0.012	1.06 (0.90–1.25)	0.468
Others	1.19 (1.01–1.40)	0.041	1.04 (0.88–1.24)	0.639

(Continued)

**Table 2** (Continued).

	Univariate Analysis HR (95% CI)	P-value	Multivariate Analysis HR (95% CI)	P-value
AFP, log <sub>10</sub> ng/dL	1.16 (1.11–1.21)	< 0.001	1.08 (1.03–1.14)	0.001
Primary anti-HCC treatment				
Supportive care	1 (reference)		1 (reference)	
Surgical resection	0.18 (0.12–0.26)	< 0.001	0.24 (0.16–0.35)	< 0.001
Liver transplantation	0.45 (0.14–1.40)	0.166	0.28 (0.09–0.89)	0.031
Local ablation therapy	0.15 (0.06–0.36)	< 0.001	0.35 (0.14–0.87)	0.024
Transarterial therapy	0.38 (0.32–0.45)	< 0.001	0.46 (0.39–0.56)	< 0.001
Chemotherapy	0.72 (0.61–0.84)	< 0.001	0.67 (0.57–0.80)	< 0.001
Radiation therapy	0.69 (0.51–0.95)	0.021	0.61 (0.44–0.84)	0.003

**Abbreviations:** AFP, alpha-fetoprotein; ALT, alanine aminotransferase; CI, confidence interval; HBV, hepatitis B virus; HCV, hepatitis C virus; HR, hazard ratio.

0.778, 0.794, and 0.784, respectively, demonstrating the discriminative ability. We then created decision trees. A schematic representation of the decision tree of each model is presented in [Supplementary Figure 2](#). Using the XGBoost model, we identified the feature importance for 3-, 6-, and 12-months survival based on the magnitude of the gain value obtained for each variable. As shown in [Figure 3](#), the primary anti-HCC treatment was a discriminative feature for predicting survival.

## Model Application

Our prediction model is available online (<https://metastatic-hcc.onrender.com/form>). For example, a male patient aged 50 years with HBV-induced, single, and 8.4 cm sized HCC with bone metastasis was part of our cohort. At the time of diagnosis, his BMI was 21.2 kg/m<sup>2</sup>, he had Child-Pugh class A, ALT of 25 IU/L, platelet count of 211x10<sup>3</sup>/mm<sup>3</sup>, sodium of 142 mmol/L, and AFP level of 16,600 ng/dL. The patient then underwent transarterial therapy. According to our model, the patient survival probabilities at 3, 6, and 12 months were 95.0%, 65.3%, and 17.9%, respectively ([Figure 4](#)). He died 9.0 months after the initial HCC diagnosis.

## Discussion

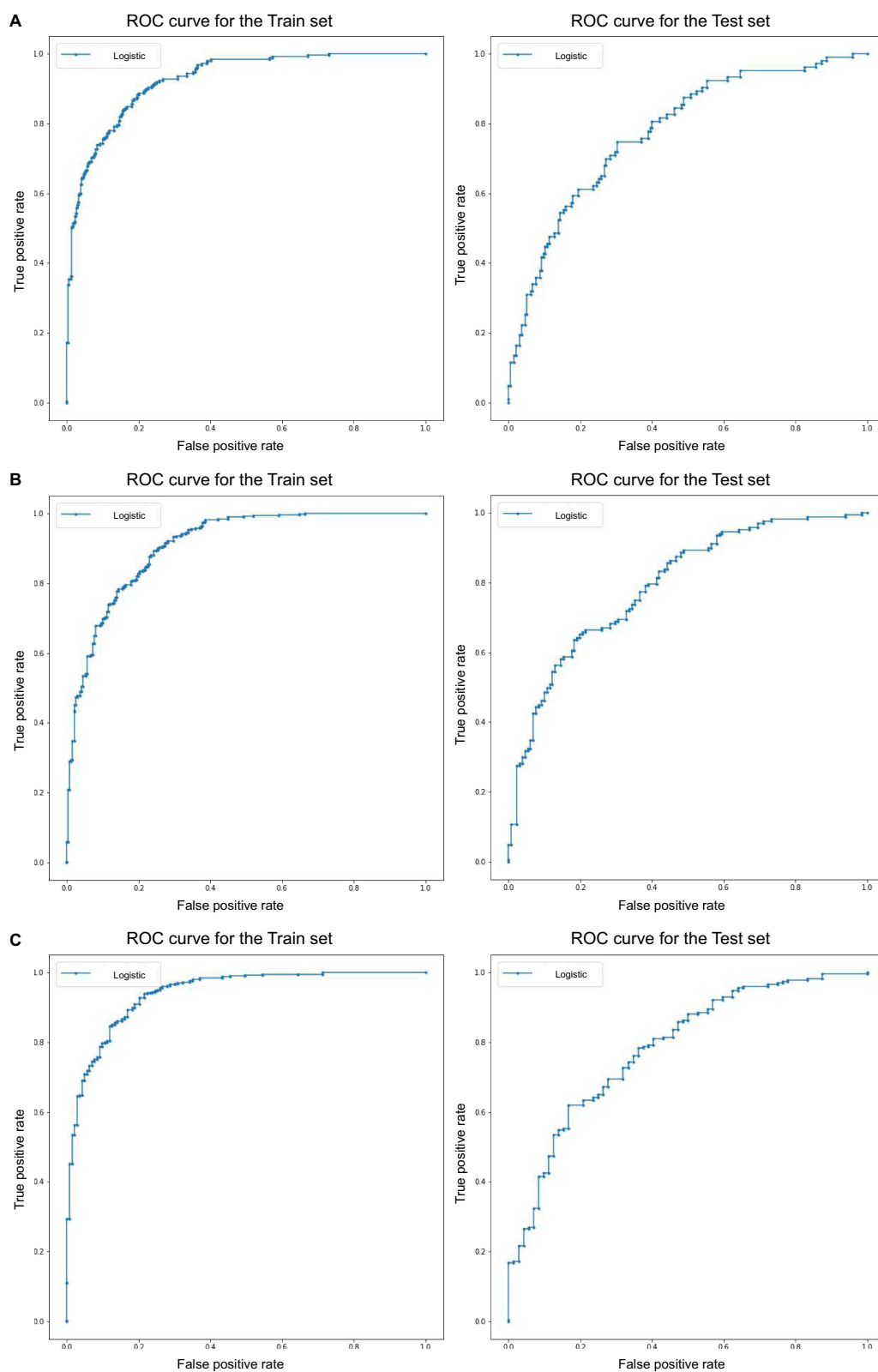
We developed a survival prediction model using XGBoost for patients with HCC and extrahepatic metastasis. We utilized a large volume of qualified data from multi-centers and used variables commonly evaluated in real-world practice. Our model achieved AUC values of 0.778, 0.794, and 0.784 for the 3-, 6-, and 12-month survival probabilities, respectively.

**Table 3** Performance Metrics of 3-, 6-, and 12-Months Survival Prediction Models

	Accuracy	Precision	Recall	F1 score	AUC
3-months					
Train set	0.849	0.840	0.706	0.767	0.925
Test set	0.735	0.654	0.495	0.564	0.778
6-months					
Train set	0.846	0.843	0.932	0.885	0.910
Test set	0.715	0.710	0.832	0.756	0.794
12-months					
Train set	0.908	0.916	0.943	0.944	0.941
Test set	0.792	0.833	0.901	0.869	0.784

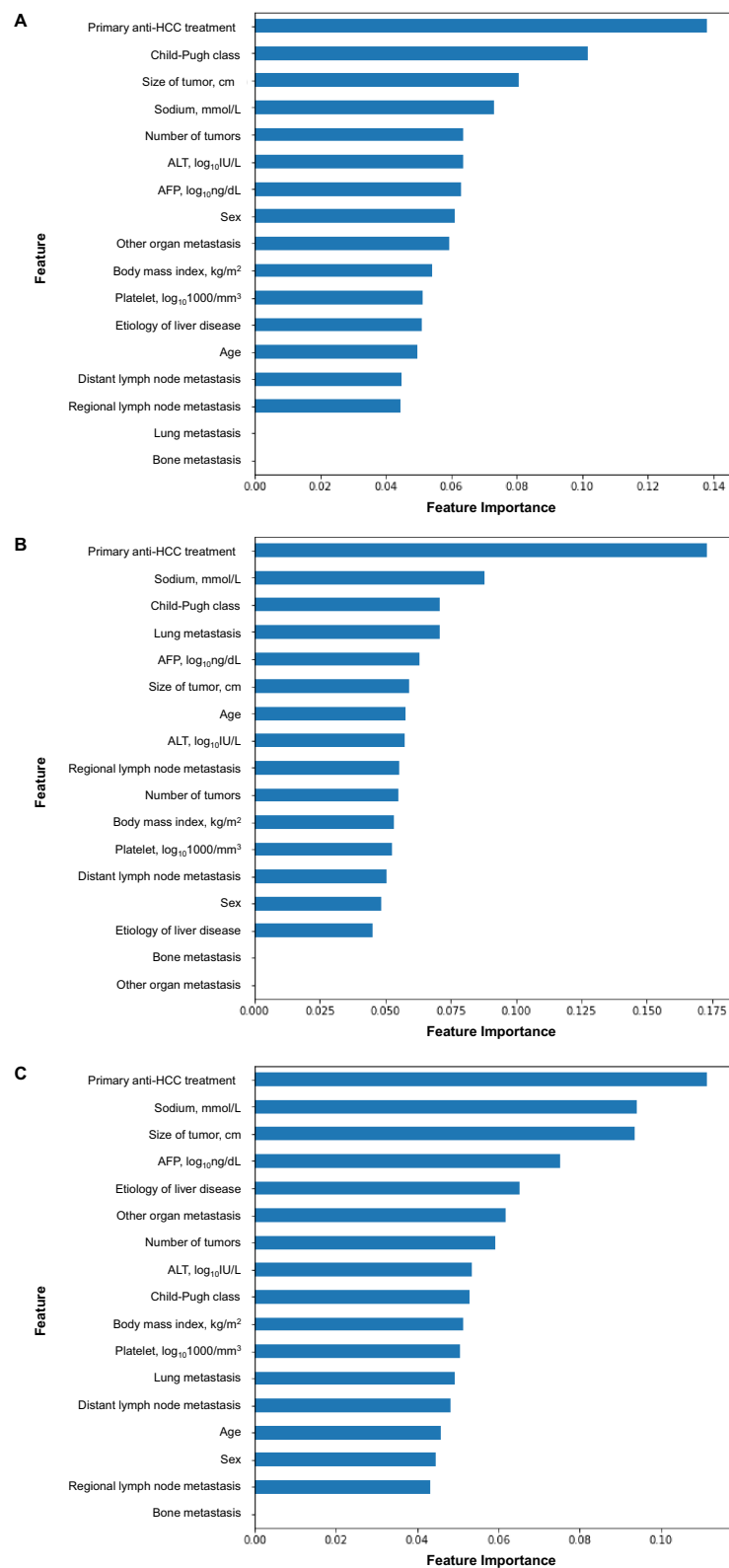
**Abbreviation:** AUC, area under the receiver operating characteristic curve.





**Figure 2** Logistic regression analysis in the test cohort. Area under the receiver operating characteristic curve for (A) 3-, (B) 6-, and (C) 12- months.





**Figure 3** Feature importance of (A) 3-, (B) 6-, and (C) 12- months survival prediction model in hepatocellular carcinoma patients with extrahepatic metastasis.

## Survival probability prediction of patients with hepatocellular carcinoma and extrahepatic metastasis

Model (months)  
3

Patients-related factors

Sex  
Male

Age  
50

BMI (kg/m<sup>2</sup>)  
21

Liver-related factors

Etiology of liver disease  
HBV infection

Child-Pugh class  
A

Platelets, 1000/mm<sup>3</sup>  
211

ALT, IU/L  
25

Sodium, mmol/L  
142

Tumor-related factors

Number of tumors  
1

Size of tumor (cm)  
7-10

Regional lymph node metastasis  
Absent

Lung metastasis  
Absent

Bone metastasis  
Present

Distant lymph node metastasis  
Absent

Other organ metastasis  
Absent

AFP, ng/dL  
16600

Primary anti-HCC treatment  
Transarterial therapy

Predict

Survival probability  
95.0%

**Figure 4** An example of applying our model to estimate survival probability at 3 months.

Using our model, patients with HCC and extrahepatic metastasis can obtain more detailed information about their prognosis and enable individualized treatment.

Being diagnosed with advanced cancer brings anxiety and uncertainty with the anticipation of suffering and fear of death.<sup>18,19</sup> The patients want to participate in their care process with necessary information for the best decisions. They hope to be prepared for their death making certain arrangements beforehand. The caregivers are also forced to face emotional and practical hardships.<sup>20</sup> They feel anxiety and depression affecting mental health before and after the

patient's death. The financial burden caused by medical care and loss of productivity often overwhelms their financial condition. Therefore, the reliable anticipation of life expectancy is crucial for patients and their loved ones.

Despite active surveillance, about 12–18% of patients with HCC are diagnosed with extrahepatic metastasis.<sup>21–24</sup> In our study, 11.6% of patients with HCC had extrahepatic metastasis. The frequency of metastatic sites differs according to the studies. A Japanese study including 151 patients with metastatic HCC reported that the lung (47.0%), regional lymph node (42.4%), and bone (37.1%) were frequent sites for HCC metastasis at 2007.<sup>23</sup> Another study conducted in the USA in 2014 using claims data revealed that the lung (30.8%), peritoneum (19.0%), bone (15.9%), and lymph nodes below the diaphragm (11.2%) were the common sites of HCC.<sup>21</sup> We found that metastasis to lung and bone had a negative effect on survival. In a European study, it was observed that the presence of lung metastasis reduced survival, whereas bone metastasis did not have a similar effect.<sup>25</sup> The association between metastatic site and survival is an intriguing finding, which could be attributed to distinct metastatic mechanisms. Lung metastasis is commonly associated with hematogenous spread through encapsulating tumor clusters and macro trabecular-massive subtypes.<sup>26–28</sup> It is expected to facilitate tumor spread through blood circulation. Whereas, lymph node metastasis occurs through increased lymphangiogenesis of the tumor along with epithelial–mesenchymal transition.<sup>29</sup> The metastatic lymph node has high immune cell infiltration with more fibrous tumor stroma than the lung reflecting immune system activation against HCC.<sup>30</sup> Lastly, bone metastasis has both hematogenous and lymphatic spread features.<sup>28</sup>

In our study, liver function and tumor burden played a critical role in prognosis consistent with previous studies.<sup>23,31,32</sup> In particular, various anti-HCC treatment strategies resulted in positive responses. Intrahepatic HCC-directed therapies such as surgery, local ablation therapy, transarterial therapy, and radiation therapy have improved survival and are recommended for patients with HCC and extrahepatic metastasis.<sup>2,3,32–36</sup> A Korean study reported that 13.1% of the patients with HCC and extrahepatic metastasis, who were treated with multimodal strategies obtained objective intrahepatic tumor response and gained improved survival compared with their counterparts.<sup>32</sup>

We employed the XGBoost algorithm for survival prediction and model development by leveraging a large patient dataset. XGBoost is a powerful machine learning algorithm that can handle high-dimensional data while reducing the risk of overfitting by automatically selecting and utilizing important factors.<sup>7</sup> Moreover, it utilizes ensemble learning techniques, which enables decision trees to improve performance and reduce learning time.<sup>7</sup> XGBoost also provides feature importance and a decision tree that enhances interpretability compared to previous machine learning algorithms.<sup>37</sup>

There are several limitations in our study. First, the median survival of our patients was only 4.0 months, which is relatively shorter than that reported in previous studies.<sup>2</sup> This shorter survival time may be attributed to the inclusion of heterogeneous patients, especially those with poorer performance or liver function who were left untreated compared to other studies. Nevertheless, including all patients is reasonable, as our goal was to build a model based on real-world data. Second, our study did not include the latest medications, such as atezolizumab-bevacizumab, durvalumab-tremelimumab, lenvatinib, cabozantinib, or ramucirumab, as our patient cohort was enrolled from 2013 to 2018.<sup>1–4</sup> Due to the inherent limitations of our data, our model may not be directly applicable to patients who have received recent chemotherapeutic agents, and it will need continuous updates, including data from patients undergoing chemotherapy. Third, it may be difficult to generalize our model to different demographic populations because of the lack of external validation. In addition, we were unable to obtain or utilize data that might play a crucial role in survival, such as ECOG performance, cirrhosis, or treatment response after the initial anti-HCC treatment. Therefore, prospective validation and improvement of our model with a more diverse patient population and precise clinical data are required to overcome these limitations.

## Conclusion

In conclusion, we successfully developed an accurate and reliable model for predicting the survival probability of HCC patients with extrahepatic metastasis using the XGBoost algorithm. Our model is easy to use and requires simple clinical data, making it accessible to both physicians and patients. We anticipate that our model will aid in individualized survival time estimation and provide valuable information for clinical decision making, ultimately leading to improved survival outcomes.

## Abbreviations

AFP, alpha-fetoprotein; ALT, alanine transaminase; AST, aspartic acid transaminase; AUC, area under the receiver operation characteristic curve; BMI, body mass index; CI, confidence interval; ECOG, eastern cooperative oncology group; FIB-4, fibrosis-4; HCC, hepatocellular carcinoma; HR, hazard ratio; MELD, model for end-stage liver disease; SD, standard deviation; TACE, transarterial chemoembolization; XGBoost, eXtreme Gradient Boosting.

## Data Sharing Statement

Data are subject to third-party restrictions. It was provided by the Korean Liver Cancer Association and the Ministry of Health and Welfare, Korea Central Cancer Registry.

## Ethics Approval Statement

This study was approved by the Institutional Review Board of Catholic University of Korea (IRB No. SC22ZISE0092).

## Consent for Publication

All authors approved the final manuscript.

## Funding

There is no funding to report.

## Disclosure

The authors declare that there are no competing interests in this work.

## References

1. Reig M, Forner A, Rimola J, et al. BCLC strategy for prognosis prediction and treatment recommendation: the 2022 update. *J Hepatol*. 2022;76(3):681–693. doi:10.1016/j.jhep.2021.11.018
2. Korean Liver Cancer Association (KLCA), National Cancer Center (NCC) Korea. 2022 KLCA-NCC Korea practice guidelines for the management of hepatocellular carcinoma. *Clin Mol Hepatol*. 2022;28(4):583–705. doi:10.3350/cmh.2022.0294
3. Omata M, Cheng AL, Kokudo N, et al. Asia-Pacific clinical practice guidelines on the management of hepatocellular carcinoma: a 2017 update. *Hepatol Int*. 2017;11(4):317–370. doi:10.1007/s12072-017-9799-9
4. European Association for the Study of the Liver. EASL Clinical Practice Guidelines: management of hepatocellular carcinoma. *J Hepatol*. 2018;69(1):182–236. doi:10.1016/j.jhep.2018.03.019
5. Dhanasekaran R, Bandoh S, Roberts LR. Molecular pathogenesis of hepatocellular carcinoma and impact of therapeutic advances. *F1000Res*. 2016;5:879. doi:10.12688/f1000research.6946.1
6. Giannini EG, Farinati F, Ciccarese F, et al. Prognosis of untreated hepatocellular carcinoma. *Hepatology*. 2015;61(1):184–190. doi:10.1002/hep.27443
7. Chen TQ, Guestrin C. XGBoost: a scalable tree boosting system. Proceedings of the 22nd Acm Sigkdd International Conference on Knowledge Discovery and Data Mining (KDD '16); 2016. doi:10.1145/2939672.2939785.
8. Budholiya K, Shrivastava SK, Sharma V. An optimized XGBoost based diagnostic system for effective prediction of heart disease. *J King Saud Univ Comput Inf Sci*. 2022;34(7):4514–4523. doi:10.1016/j.jksuci.2020.10.013
9. Kim JOR, Jeong YS, Kim JH, Lee JW, Park D, Kim HS. Machine learning-based cardiovascular disease prediction model: a cohort study on the Korean National health insurance service health screening database. *Diagnostics*. 2021;11(6):943. doi:10.3390/diagnostics11060943
10. Luo XQ, Yan P, Duan SB, et al. Development and validation of machine learning models for real-time mortality prediction in critically ill patients with sepsis-associated acute kidney injury. *Front Med*. 2022;9:853102. doi:10.3389/fmed.2022.853102
11. Zhang Y, Zhang X, Razbek J, et al. Opening the black box: interpretable machine learning for predictor finding of metabolic syndrome. *BMC Endocr Disord*. 2022;22(1):214. doi:10.1186/s12902-022-01121-4
12. Carmona P, Dwekat A, Mardawi Z. No more black boxes! Explaining the predictions of a machine learning XGBoost classifier algorithm in business failure. *Res Int Bus Finance*. 2022;61:101649. doi:10.1016/j.ribaf.2022.101649
13. Xu XL, Jiang LS, Wu CS, et al. The role of fibrosis index FIB-4 in predicting liver fibrosis stage and clinical prognosis: a diagnostic or screening tool? *J Formos Med Assoc*. 2022;121(2):454–466. doi:10.1016/j.jfma.2021.07.013
14. Cheema JR. Some general guidelines for choosing missing data handling methods in educational research. *J Mod Appl Stat Methods*. 2014;13(2):53–75. doi:10.22237/jmasm/1414814520
15. Ditzler G, LaBarck J, Ritchie J, Rosen G, Polikar R. Extensions to online feature selection using bagging and boosting. *IEEE Trans Neural Netw Learn Syst*. 2018;29(9):4504–4509. doi:10.1109/TNNLS.2017.2746107
16. Fawcett T. An introduction to ROC analysis. *Pattern Recognit Lett*. 2006;27(8):861–874. doi:10.1016/j.patrec.2005.10.010
17. Render [homepage on the Internet]. The fastest way to host all your web apps; 2023. Available from: <https://render.com/>. Accessed October 29, 2023.

18. von Blanckenburg P, Riera Knorrenschild J, Hofmann M, et al. Expectations, end-of-life fears and end-of-life communication among palliative patients with cancer and caregivers: a cross-sectional study. *BMJ Open*. 2022;12(5):e058531. doi:10.1136/bmjopen-2021-058531
19. Mazzocco K, Masiero M, Carriero MC, Pravettoni G. The role of emotions in cancer patients' decision-making. *Ecancermedicalscience*. 2019;13:914. doi:10.3332/ecancer.2019.914
20. Jacobs JM, Shaffer KM, Nipp RD, et al. Distress is interdependent in patients and caregivers with newly diagnosed incurable cancers. *Ann Behav Med*. 2017;51(4):519–531. doi:10.1007/s12160-017-9875-3
21. Abbas A, Medvedev S, Shores N, et al. Epidemiology of metastatic hepatocellular carcinoma, a nationwide perspective. *Dig Dis Sci*. 2014;59(11):2813–2820. doi:10.1007/s10620-014-3229-9
22. Katyal S, Oliver JH, Peterson MS, Ferris JV, Carr BS, Baron RL. Extrahepatic metastases of hepatocellular carcinoma. *Radiology*. 2000;216(3):698–703. doi:10.1148/radiology.216.3.r00se24698
23. Uka K, Aikata H, Takaki S, et al. Clinical features and prognosis of patients with extrahepatic metastases from hepatocellular carcinoma. *World J Gastroenterol*. 2007;13(3):414–420. doi:10.3748/wjg.v13.i3.414
24. Feng J, He Y, Wan J, Chen Z. Pulmonary metastases in newly diagnosed hepatocellular carcinoma: a population-based retrospective study. *HPB (Oxford)*. 2020;22(9):1295–1304. doi:10.1016/j.hpb.2019.12.004
25. Schütte K, Schinner R, Fabritius MP, et al. Impact of extrahepatic metastases on overall survival in patients with advanced liver dominant hepatocellular carcinoma: a subanalysis of the SORAMIC trial. *Liver Cancer*. 2020;9(6):771–786. doi:10.1159/000510798
26. Sneag DB, Krajewski K, Giardino A, et al. Extrahepatic spread of hepatocellular carcinoma: spectrum of imaging findings. *AJR Am J Roentgenol*. 2011;197(4):W658–64. doi:10.2214/AJR.10.6402
27. Lin YL, Li Y. Study on the hepatocellular carcinoma model with metastasis. *Genes Dis*. 2020;7(3):336–350. doi:10.1016/j.gendis.2019.12.008
28. Yuan X, Zhuang M, Zhu X, et al. Emerging perspectives of bone metastasis in hepatocellular carcinoma. Systematic Review. *Front Oncol*. 2022;12:943866. doi:10.3389/fonc.2022.943866
29. Roy S, Banerjee P, Ekser B, et al. Targeting lymphangiogenesis and lymph node metastasis in liver cancer. *Am J Pathol*. 2021;191(12):2052–2063. doi:10.1016/j.ajpath.2021.08.011
30. Woo HY, Rhee H, Yoo JE, et al. Lung and lymph node metastases from hepatocellular carcinoma: comparison of pathological aspects. *Liver Int*. 2022;42(1):199–209. doi:10.1111/liv.15051
31. Uchino K, Tateishi R, Shiina S, et al. Hepatocellular carcinoma with extrahepatic metastasis: clinical features and prognostic factors. *Cancer*. 2011;117(19):4475–4483. doi:10.1002/cncr.25960
32. Jung SM, Jang JW, You CR, et al. Role of intrahepatic tumor control in the prognosis of patients with hepatocellular carcinoma and extrahepatic metastases. *J Gastroenterol Hepatol*. 2012;27(4):684–689. doi:10.1111/j.1440-1746.2011.06917.x
33. Komatsu S, Kido M, Tanaka M, et al. Clinical significance of hepatectomy for hepatocellular carcinoma associated with extrahepatic metastases. *Dig Surg*. 2020;37(5):411–419. doi:10.1159/000507436
34. Chang WI, Kim BH, Kim YJ, Yoon JH, Jung YJ, Chie EK. Role of radiotherapy in Barcelona clinic liver cancer stage C hepatocellular carcinoma treated with sorafenib. *J Gastroenterol Hepatol*. 2022;37(2):387–394. doi:10.1111/jgh.15722
35. Shao YY, Wang SY, Lin SM. Management consensus guideline for hepatocellular carcinoma: 2020 update on surveillance, diagnosis, and systemic treatment by the Taiwan liver cancer association and the gastroenterological society of Taiwan. *J Formos Med Assoc*. 2021;120(4):1051–1060. doi:10.1016/j.jfma.2020.10.031
36. Xie DY, Ren ZG, Zhou J, Fan J, Gao Q. 2019 Chinese clinical guidelines for the management of hepatocellular carcinoma: updates and insights. *Hepatobiliary Surg Nutr*. 2020;9(4):452–463. doi:10.21037/hbsn-20-480
37. Li Z. Extracting spatial effects from machine learning model using local interpretation method: an example of SHAP and XGBoost. *Comput Environ Urban Syst*. 2022;96:101845. doi:10.1016/j.compenurbysys.2022.101845