

Original Article

A predictive model for the diagnosis of non-alcoholic fatty liver disease based on an integrated machine learning method

Xuefeng Ma¹, Chao Yang², Kun Liang², Baokai Sun¹, Wenwen Jin¹, Lizhen Chen¹, Mengzhen Dong¹, Shousheng Liu³, Yongning Xin¹, Likun Zhuang³

¹Department of Infectious Disease, Qingdao Municipal Hospital, Qingdao University, Qingdao 266000, Shandong, China; ²Department of Infectious Disease, The Affiliated Hospital of Qingdao University, Qingdao 266000, Shandong, China; ³Clinical Research Center, Qingdao Municipal Hospital, Qingdao University, Qingdao 266000, Shandong, China

Received May 11, 2021; Accepted October 12, 2021; Epub November 15, 2021; Published November 30, 2021

Abstract: Diagnostic markers for non-alcoholic fatty liver disease (NAFLD) are still needed for screening individuals at risk. In recent years, the machine learning method was used to search for the diagnostic markers of multiple diseases. In this study, we developed and validated a machine learning model to diagnose NAFLD using laboratory indicators. NAFLD patients and non-NAFLD controls were recruited in the training and validation cohorts. The laboratory indicators of the participants in the training cohort were collected, and six indicators including alanine aminotransferase/aspartate aminotransferase (ALT/AST), white blood cells (WBC), alpha-L-fucosidase (AFU), hemoglobin (Hb), triglycerides (TG) and gamma-glutamyl transpeptidase (GGT) were screened out with higher weights by an integrate machine learning method. The areas under the receiver operating characteristic curves (AUROCs) for the selected indicators using logistic regression (LR), random forest (RF) and support vector machine (SVM) were 0.814, 0.837 and 0.810, respectively. Then the binary logistic regression was used to construct the predictive model. What's more, the AUROC of the predicted model was 0.732 in the validation cohort of patients with NAFLD. And the combined AUROC of the six parameters was 0.716 in the mouse model fed with high-fat diet (HFD). In summary, we created a predictive model with six laboratory indicators for the diagnosis of NAFLD based on the machine learning method, which has the potential value for the diagnosis of the NAFLD.

Keywords: NAFLD, diagnosis, machine learning, laboratory indicator

Introduction

Non-alcoholic fatty liver disease (NAFLD) is one of the most common chronic liver diseases globally with an estimated prevalence of about 24% in North America and about 30% in Asia [1, 2]. NAFLD is marked by excessive intrahepatic fat deposition [3]. The disease spectrum of NAFLD ranges from simple steatosis to non-alcoholic steatohepatitis (NASH), fibrosis, cirrhosis and even hepatocellular carcinoma (HCC) [4], which could bring heavy burdens on the health care system [5, 6].

Although liver pathology was the gold standard for the diagnosis of NAFLD, it is invasive, which limits the wide application. Other diagnostic

methods for NAFLD were based on imaging examinations such as B ultrasound, computed tomography (CT) or Fibroscan according to the Practice Guidance of the American Association for the Study of Liver Diseases (AASLD) Practice Guidelines [7]. However, CT and Fibroscan are not convenient for patients, and B ultrasound was limited in relatively low accuracy and specificity. Recently, blood indexes, which are convenient, low-cost and readily available, are ideal tools for the diagnosis of diseases. Several prediction models based on blood indexes have been constructed for the diagnosis of NAFLD. Elevated alanine aminotransferase (ALT) level is the predominant finding for the diagnosis of NAFLD, but elevated ALT level also occurred in other liver diseases such as

A predictive model for NAFLD

hepatitis B virus (HBV) infection and intrahepatic cholangitis [8]. SteatoTest for the diagnosis of NAFLD is a logistic regression model containing 12 predicting indicators including a2-macroglobulin (A2M), apolipoprotein A1 (Apo A1), haptoglobin, gamma-glutamyl transpeptidase (GGT) levels, total bilirubin, cholesterol, triglycerides, glucose, age, gender, and body mass index [9], while it has not been proved as a practical model based on the combination of blood indexes with high sensitivity and specificity.

In recent years, machine learning methods including filter method, wrapper method and the embedded method were developed to select predicting parameters among all available indicators with maximum data and minimum bias to predict diseases. Machine learning could reveal the complex relationships among the indicators, which was useful for the diagnosis of diseases. Li et al. built multiparametric ultrasomics which could improve discrimination of significant fibrosis in chronic hepatitis B patients compared with mono or dual modalities [10]. What's more, Liu et al. built an artificial neural network model that was useful for evaluating the probability of progression-free survival in patients with HCC [11].

In this study, we aimed to develop a predictive model for diagnosing NAFLD depending on the laboratory parameters using an integrated machine learning method and validated the diagnostic effects of the model for diagnosing NAFLD at the population and animal levels.

Methods

Research subjects

Inclusion criteria for NAFLD patients and non-NAFLD patients in this study: (1) 18-65 years of age; (2) Patients were negative for hepatitis B surface antigen, hepatitis C virus-RNA and hepatitis B virus DNA; (3) No clinical symptoms or signs of infection, no liver disorders or other critical diseases, and no fractures, osteoporosis, or tumors; (4) No pregnancy for women; (5) No drinks or no more than 70 g ethanol per week for women (about one standard drink daily), and 140 g for men (two standard drinks daily).

Exclusion criteria for NAFLD patients and non-NAFLD patients in this study: (1) Patients with

an active implantable medical device (such as pacemaker or defibrillator); (2) Hematological diseases or diseases that may influence the parameters of blood cell counts; (3) Patients who had undergone liver transplantation, patients with cardiac failure and/or significant valvular disease.

Definition of NAFLD: Patients with and without NAFLD were first distinguished by pathology. If the pathology was not available, B ultrasound or CT was used. The diagnostic criteria of NAFLD were followed by the practice guidance from the American Association [12].

According to the criteria above, a total of 45 NAFLD patients and 53 non-NAFLD controls from Qingdao Municipal Hospital were recruited in the training group. Informed consent was signed by every participant. This study was approved by the Ethics Committee of the Qingdao Municipal Hospital (2019Y006). 81 NAFLD patients and 87 non-NAFLD controls from the Affiliated Hospital of Qingdao University were involved in the validation cohort. Informed consent was signed by every participant. This study was also approved by the Ethics Committee of the Affiliated Hospital of Qingdao University (QYFYWZLL26473).

Clinical data collection

The data for the clinical laboratory indicators of the first admission were collected from the Qingdao Municipal Hospital and the Affiliated Hospital of Qingdao University ([Supplementary Tables 1](#) and [2](#)). The data for the imaging examinations were also collected at the meantime.

Animal experiments

Eight-week-old male C57BL/6J mice were fed with high-fat diet (HFD) or chow diet (CD). Twelve weeks later, all mice were sacrificed. The hematoxylin-eosin (HE) staining of liver tissues was used to verify the successful establishment of the NAFLD mice model. WBC was calculated by XFA6030 Animal Blood Cell Analyzer (prolong, Beijing, China). ALT, aspartate aminotransferase (AST), triglycerides (TG), hemoglobin (Hb), and alpha-L-fucosidase (AFU) were tested by enzyme linked immunosorbent assay kit (Shanghai Enzyme-linked Biotechnology Co., Ltd., Shanghai, China). Animal experiments were approved by the Animal Experi-

A predictive model for NAFLD

Table 1. Partial differential characteristics of subjects with or without NAFLD in the training cohort

Characters	NAFLD n=53	No NAFLD n=45	P value
Gender (male/female, n)	32/21	25/20	0.02
ALT (U/L)	31.21±21.54	19.30±12.93	<0.001
ALT to AST	1.33±0.64	0.90±0.29	<0.001
GGT (U/L)	35.88±29.93	21.95±13.96	0.02
AFU (U/L)	26.45±7.91	23.32±7.12	0.046
WBC (× 10 ⁹ /L)	6.40±1.70	5.30±1.31	<0.001
TG (mmol/L)	1.83±1.26	1.25±0.53	<0.001
Hb (g/L)	144.98±11.69	138.22±15.76	0.03
NEUT (10 ⁹ /L)	3.40±1.11	2.90±0.94	0.02
LYM (10 ⁹ /L)	2.78±2.52	2.00±0.62	0.03
MON (10 ⁹ /L)	0.40±0.12	0.34±0.12	0.02
BASO (10 ⁹ /L)	0.04±0.02	0.03±0.02	0.04
UA (μmol/L)	360.34±86.08	312.26±70.99	<0.001
HDL (mmol/L)	1.24±0.28	1.39±0.31	0.02
APO A1/APO B	1.35±0.39	1.62±0.70	0.04

mental Ethical Committee of Qingdao University (AHQU-MAL20180504-1).

Statistical analysis

The weight values of candidate biomarkers were calculated by the method of integrated machine learning (Applied Protein Technology, Shanghai, China) [13-15]. To validate the effects of selected candidate biomarkers in classification models, three commercialized machine learning models including logistic regression (LR) [16], random forest (RF) [17] and support vector machine (SVM) [18] were used respectively. The diagnostic values of the indicators were evaluated by the receiver operating characteristic (ROC) curve and the areas under the receiver operating characteristic curves (AUROC). Pearson correlation coefficient (r) was calculated between the candidate biomarkers. When the value of r was more than 0.6, the correlation was defined as strong [19]. Data were analyzed by R version 3.6.1 (R Foundation for Statistical Computing, Vienna, Austria). Data were expressed as mean ± standard deviation (SD). Missing values were imputed using mean imputation. The difference between subgroups was analyzed by chi-square test for categorical parameters and Student's t-test was used for continuous parameters.

Results

Patient characteristics

In the training cohort including 53 NAFLD patients and 45 non-NAFLD controls, the percentage of male was 60.38% in the NAFLD group, which was higher than that in the non-NAFLD controls (P<0.05). There were also significant differences in ALT/AST, ALT, white blood cells (WBC), AFU, Hb, TG, GGT, neutrophilic granulocyte count (NEUT), lymphocyte count (LYM), monocytes count (MON), basophils (BASO), uric acid (UA), high-density lipoprotein (HDL) and apolipoprotein A1/apolipoprotein B (APO A1/APO B) between the two groups (**Table 1**).

Six parameters were selected using an integrated machine learning

method

For the training cohort, the weight value of each biomarker was calculated by an integrated machine learning method. The higher value of the weight means the greater contribution of the biomarker in distinguishing the patients with NAFLD from the non-NAFLD controls. The top six candidate biomarkers were WBC, ALT to AST, AFU, Hb, TG and GGT (**Figure 1A**).

The cumulative AUROC chart was used to evaluate the impact of combined biomarkers and the sequence for the combination of indicators were determined according to the weight values. As shown in **Figure 1B**, the AUROC of the top six indicators was 0.772 and when the 7th indicator was added, the AUROC of the combined biomarkers was suddenly reduced. Consequently, we selected six parameters for further analysis.

Verification of the candidate biomarkers

ROC curves of three models including LR, RF and SVM demonstrated that the selected candidate biomarkers have excellent effects for classification, and the AUROC values of LR, RF, and SVM were 0.814, 0.837 and 0.810, respectively (**Figure 2A**). The importance coefficients of the candidate biomarkers were calculated by RF to compare the contribution of

A predictive model for NAFLD

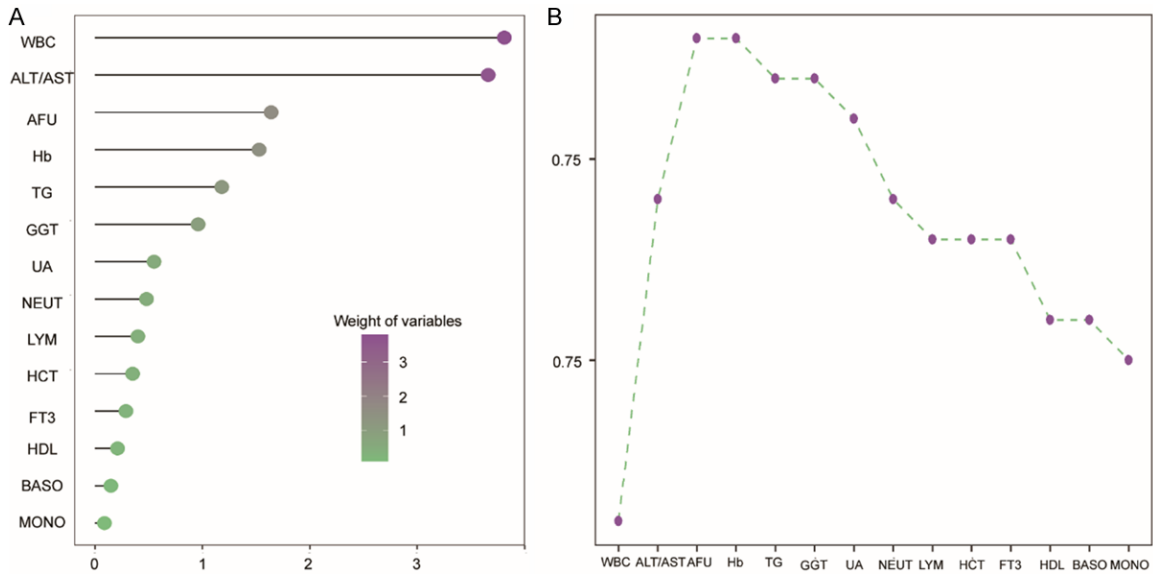
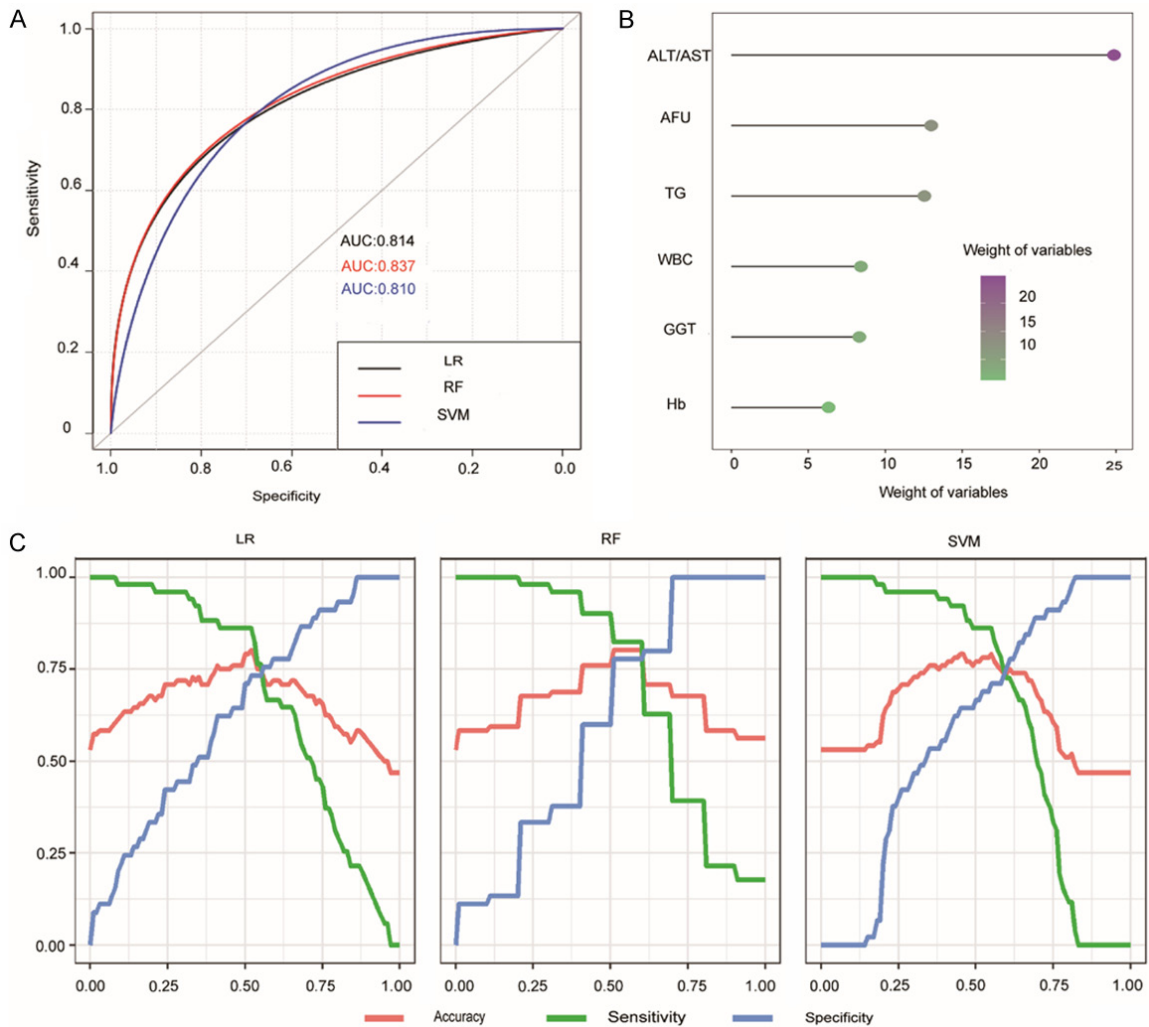


Figure 1. Six parameters were selected using an integrated machine learning method in the training cohort of patients with NAFLD. A. Weight values of the variables; B. The cumulative AUC chart of the candidate biomarkers.



A predictive model for NAFLD

Figure 2. Verification of the candidate biomarkers in the training cohort of patients with NAFLD. A: The ROC curves of the candidate biomarkers for the diagnosis of NAFLD using LR, SVM and RF methods; B: The importance of the candidate biomarkers calculated by RF; C: Curves of accuracy, sensitivity, and specificity of the candidate biomarkers for the diagnosis of NAFLD evaluated by LR, RF and SVM.

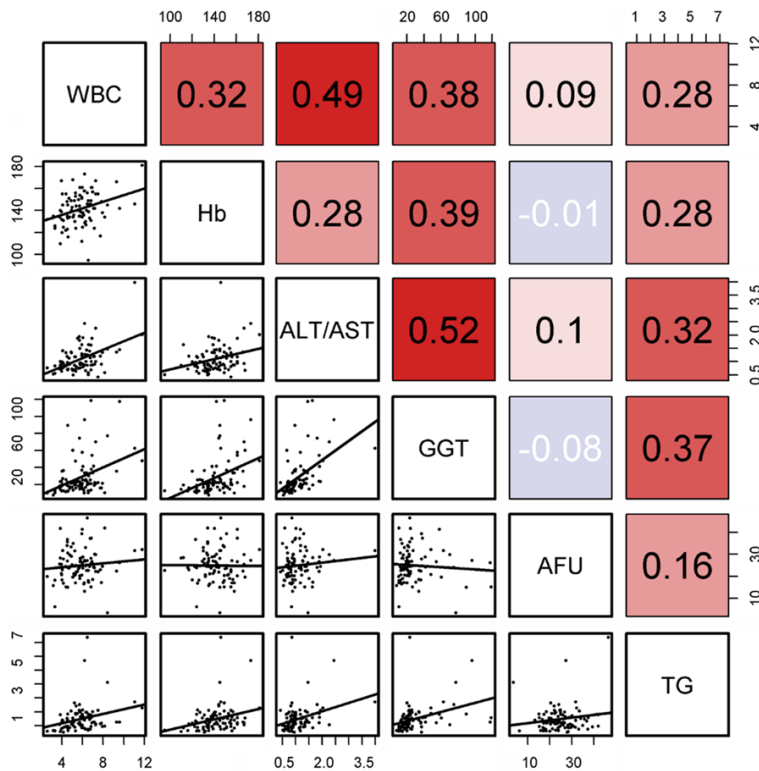


Figure 3. Correlations of the candidate biomarkers in the training cohort of patients with NAFLD. Pearson correlation coefficient (r) was used to evaluate the correlations.

each biomarker in the model (**Figure 2B**). Accuracy, sensitivity, and specificity evaluated by LR, RF and SVM were revealed in **Figure 2C**. The results of verification demonstrated that the selected biomarkers have excellent effects of classification for NAFLD patients in the training cohort.

Correlation of the candidate biomarkers

We conducted the correlation analysis among the selected biomarkers. If there was a strong correlation between the two biomarkers, one of them would be eliminated. As shown in **Figure 3**, no strong correlations between the selected biomarkers were observed, and none of the biomarkers in this study was eliminated.

Construction of the predicted model

The diagnostic panel of candidate biomarkers for NAFLD was built by logistic regression algorithm. The logical regression coefficients of

WBC, ALT to AST, AFU, Hb, TG and GGT were 0.29, 2.23, 0.07, 0.01, 0.72 and 0.01, respectively (**Supplementary Table 3**). The equation was illustrated as following:

$$\text{Risk score} = \exp(0.29 \times \text{WBC} + 2.23 \times \frac{\text{ALT}}{\text{AST}} + 0.07 \times \text{AFU} + 0.01 \times \text{Hb} + 0.72 \times \text{TG} + 0.01 \times \text{GGT} - 8.77) / [1 + \exp(0.29 \times \text{WBC} + 2.23 \times \frac{\text{ALT}}{\text{AST}} + 0.07 \times \text{AFU} + 0.01 \times \text{Hb} + 0.72 \times \text{TG} + 0.01 \times \text{GGT} - 8.77)]$$

The best cutoff value of the risk score was 0.53 calculated by using Yoden Index (**Supplementary Table 4**). When the risk score was 0.53, the AUROC value, specificity and sensitivity for the predicted model were 0.821, 0.733 and 0.765, respectively. If the risk score of the equation was higher than the best cutoff value, the result would be positive.

The diagnostic evaluation of the predicted model in the validation cohort

81 NAFLD patients and 87 non-NAFLD controls were involved in the validation cohort from the Affiliated Hospital of Qingdao University. There were significant differences in the values of ALT, AST, ALT to AST, GGT, AFU, WBC, TG, LYM, MON, MON% BASO%, red blood cell volume distribution width (RDW), uric acid (UA), total cholesterol (TC), HDL, APO A1, APO B and APO A1/APO B between the two groups (**Table 2**). The predicted model kept its diagnostic efficacy in the validation group with the AUROC 0.732 (**Figure 4**).

Validation of the selected parameters using the mouse model

Mice fed with HFD or CD were involved in animal experiments. The results of H&E and Oil Red O staining of liver tissue sections verified

A predictive model for NAFLD

Table 2. Partial differential characteristics of subjects with or without NAFLD in the validation cohort

Characters	NAFLD group n=81	Control group n=87	P value
ALT (U/L)	31.45±19.24	23.27±15.86	<0.001
AST (U/L)	21.82±8.52	19.22±6.69	0.03
ALT to AST	1.41±0.44	1.16±0.52	<0.001
GGT (U/L)	33.78±14.53	18.20±6.47	0.02
AFU (U/L)	29.35±7.88	26.31±7.05	0.02
WBC (× 10 ⁹ /L)	6.46±1.69	5.80±1.25	0.01
TG (mmol/L)	2.03±1.89	1.39±1.14	0.02
Hb (g/L)	143.54±14.10	138.80±17.20	0.08
LYM (10 ⁹ /L)	2.19±0.75	1.91±0.59	0.01
MON (10 ⁹ /L)	0.46±0.15	0.37±0.13	<0.001
MON%	7.32±2.02	6.42±2.02	<0.001
BASO%	0.53±0.25	0.44±0.28	0.04
RDW%	12.37±0.55	22.65±13.79	<0.001
UA (μmol/L)	366.70±78.65	306.36±80.34	<0.001
TC (mmol/L)	4.54±0.54	4.15±1.60	<0.001
HDL (mmol/L)	1.33±0.30	1.46±0.36	0.04
APO A1 (g/L)	1.38±0.24	1.52±0.22	0.01
APO B (g/L)	1.00±0.24	0.85±0.19	0.01
APO A1/APO B	1.48±0.51	1.88±0.53	<0.001

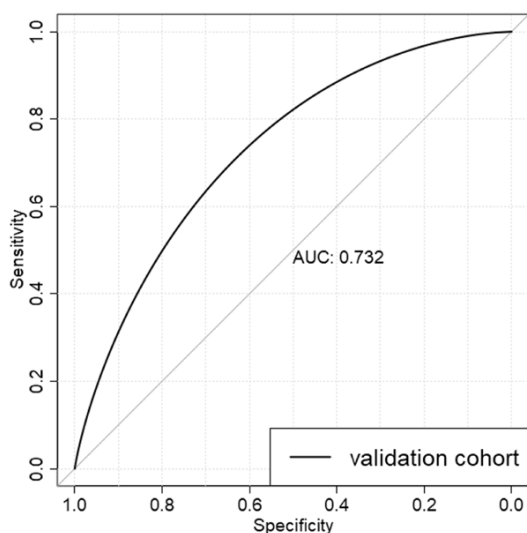


Figure 4. The ROC curve of the candidate biomarkers for the diagnosis of NAFLD in the validation cohort of NAFLD patients.

the construction of NAFLD mouse model (**Figure 5A**). Although there were no obvious differences in ALT, AST, ALT to AST, AFU, WBC, TG, Hb and GGT between the two groups (all $P > 0.05$), the average values of all the selected indicators in HFD group were higher than those

in CD group (**Supplementary Table 5**). Furthermore, the cumulative AU-ROC of the candidate biomarkers was up to 0.716 (**Figure 5B**).

Discussion

In this study, our predictive model for NAFLD, which consisted of the common laboratory indicators in hospital, showed a brilliant performance. This study aimed to screen out the combination of laboratory indicators with higher sensitivity and specificity for clinical screening of NAFLD. Although there were no relevant mechanisms revealed in this study, the indicators screened out in this study were all classic indicators and there have been many reports on the relevant mechanisms for these indicators in NAFLD.

Both ALT and AST were mainly expressed in liver cells and their levels in plasma could indicate the

damage and death of liver cells. When the liver injury occurred, ALT and AST are released from liver cells into the blood, leading to the increased serum ALT and AST levels [20]. Nanji et al. demonstrated that there was a significant correlation between the ALT/AST ratio and the degree of fatty accumulation of hepatic cells [21]. Long et al. also demonstrated that the ALT/AST ratio predicted hepatic steatosis better than either ALT or AST alone [22].

NAFLD was considered to be the liver manifestation of metabolic symptoms (MS) [23]. The relationship between WBC count and MS components had been demonstrated in some studies [24, 25]. Moreover, WBC count was often used to evaluate inflammatory status [26] and inflammation plays a significant role in the development of NAFLD [3]. In view of the points above, WBC might reflect the occurrence and progression of NAFLD. Many studies have focused on the association between WBC count and the occurrence of NAFLD [27, 28], and a previous study has clearly showed that the WBC count was a significant factor associated with NAFLD occurrence [29].

AFU is a lysosome enzyme expressed in all mammalian cells and hydrolyzes sugars con-

A predictive model for NAFLD

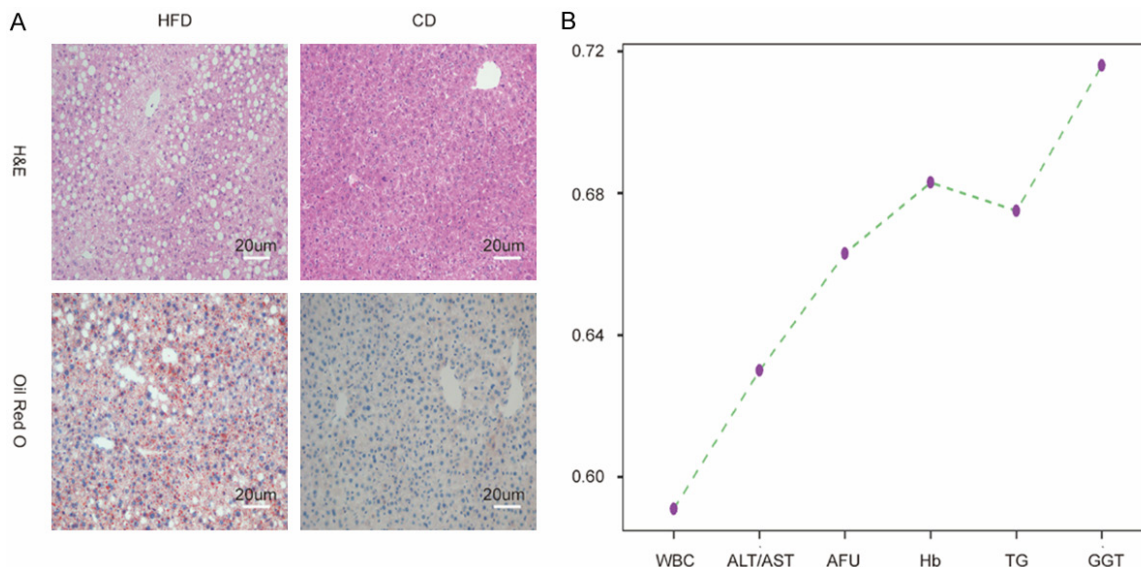


Figure 5. The Verification of the candidate parameters using the mouse model. A: H&E and Oil Red O stainings of representative liver sections in mice fed with HFD or CD. Scale bar: µm; B: The AUC Cumulative Chart of the selected parameters in the mouse model.

taining L-fucose, and AFU levels are closely associated with the occurrence of HCC [30, 31]. Blood AFU usually comes from lysosome leakage. Lipid peroxidation of liver cells modifies the functional characteristics of the cell membranes and membranes of intracellular organelles such as mitochondria and lysosomes, which might result in leakage of lysosome and the release of AFU [32, 33]. As a result, AFU may have a close relationship with NAFLD. A previous study had attempted to explore the relationship between the AFU level and NAFLD occurrence. Lu et al. suggested that AFU levels were positively associated with NAFLD occurrence and might act as an independent risk factor for NAFLD. However, the AUROC of the only AFU levels for NAFLD diagnosis was only 0.606 [34].

Hb is an iron-containing metalloprotein. A previous study revealed that the iron depletion could increase the glucose uptake and insulin signaling in hepatic cells and improve liver function in NAFLD patients [35]. Bai et al. revealed that adults with high Hb levels (14.4 µg/dl for male and 13.2 µg/dl for female) were at the greatest risk for NAFLD [36]. Chung et al. demonstrated that serum Hb level was independently associated with the risk of developing incidental metabolic syndrome or NAFLD in men [37]. Recently, Giorgio et al. indicated that elevated Hb Level had obvious relation-

ship with fibrosis in biopsy-diagnosed pediatric NAFLD patients [38]. Consistently, in a Mexican population study, there was an independent relationship between the serum Hb level and the steatosis severity [39].

GGT, which is a transmembrane protein generated in the microsomes, could play an important role in maintaining the metabolism of glutathione and act as one of the most important antioxidants in human cells [40]. Meanwhile, oxidative stress could play an essential role in the development of NAFLD [40, 41]. Jarčuška et al. found that about half of patients with NASH had elevated levels of GGT, and there was an obvious relationship between GGT and individual metabolic syndrome [42]. Furthermore, a previous study revealed that the GGT-to-platelet ratio (GPR) is better than aspartate transaminase-to-platelet ratio index and fibrosis index based on four factors (FIB-4) for diagnosing fibrosis and cirrhosis in NAFLD patients [43].

For patients with overnutrition and obesity, there is often a change of hepatic fatty acid metabolism, which could lead to the accumulation of triglycerides in hepatocytes and sometimes cause the occurrence of NAFLD [44].

Each indicator in this study has a certain tendency and limitation on the diagnosis of

A predictive model for NAFLD

NAFLD, while the predictive model using the integrated machine learning method can overcome the limitations of a single indicator and improve the diagnostic ability. Some indicators which were not common in clinical laboratory were also reported to be the potential predictors for NAFLD. Kimura et al. demonstrated that serum thrombospondin 2 level was considered as a predictor of histological activity of NAFLD [45]. Mele et al. displayed that angiopoietin-like 8 has direct relationship with the presence and severity of NAFLD in patients with Prader-Willi Syndrome [46]. In further investigations, the uncommon indicators in clinical laboratory should be considered for the construction of predictive model.

In summary, we created a predictive model with laboratory indicators for the diagnosis of NAFLD using the machine learning method, which had the potential for the diagnosis of the NAFLD and provided the basis for the predictive model of the combination of indicators.

Acknowledgements

This work was supported by grants from the National Natural Science Foundation of China [grant number: 31770837] and the Key Research and Development Program of Shandong Province [grant number: 2019GSF108148].

Disclosure of conflict of interest

None.

Abbreviations

A/G, the ratio of ALB and GLOB; APO A1, apolipoprotein A1; APO B, apolipoprotein B; AST, aspartate aminotransferase; AUROC, the areas under the receiver-operating characteristic curve; BASO, basophils; BUN, blood urea nitrogen; CK, creatine kinase; CK-MB, creatine kinase isoenzyme-MB; Cre, creatinine; DB, direct bilirubin; EOS, eosinophils count; GLOB, globulin; GLU, blood glucose; HCT, hematocrit; HDL, high-density lipoprotein; IB, indirect bilirubin; LDH, lactate dehydrogenase; LDL, low-density lipoprotein; LP, lipoproteins; LR, logistic regression; LYM, lymphocyte count; MCH, mean corpuscular hemoglobin; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume; MON, monocytes count; MPV, mean platelet volume; NAFLD, nonalcoholic fatty liver disease; NEUT, the neutrophilic

granulocyte count; PCT, platelet hematocrit; PDW, platelet distribution width; P-LCR, platelet-larger cell ratio; PLT, platelet count; RBC, red blood cell; RDW, red blood cell volume distribution width; RF, random forest; SVM, support vector machine; TB, total bilirubin; TBA, total bile acid; TC, total cholesterol; TP, total protein; UA, uric acid; WBC, white blood cell count; HBDH, alpha-hydroxybutyrate dehydrogenase; GGT, gamma-glutamyl transpeptidase.

Address correspondence to: Likun Zhuang, Clinical Research Center, Qingdao Municipal Hospital, Qingdao University, No. 5 Middle Donghai Road, Qingdao 266000, Shandong, China. Tel: +86-185-62529612; E-mail: zlk0823@163.com; Yongning Xin, Department of Infectious Disease, Qingdao Municipal Hospital, Qingdao University, No. 5 Middle Donghai Road, Qingdao 266000, Shandong, China. Tel: +86-532-82789463; E-mail: xinyongning9812@163.com

References

- [1] Li J, Zou B, Yeo YH, Feng Y, Xie X, Lee DH, Fujii H, Wu Y, Kam LY, Ji F, Li X, Chien N, Wei M, Ogawa E, Zhao C, Wu X, Stave CD, Henry L, Barnett S, Takahashi H, Furusyo N, Eguchi Y, Hsu YC, Lee TY, Ren W, Qin C, Jun DW, Toyoda H, Wong VW, Cheung R, Zhu Q and Nguyen MH. Prevalence, incidence, and outcome of non-alcoholic fatty liver disease in Asia, 1999-2019: a systematic review and meta-analysis. *Lancet Gastroenterol Hepatol* 2019; 4: 389-398.
- [2] Younossi Z, Tacke F, Arrese M, Chander Sharma B, Mostafa I, Bugianesi E, Wai-Sun Wong V, Yilmaz Y, George J, Fan J and Vos MB. Global perspectives on nonalcoholic fatty liver disease and nonalcoholic steatohepatitis. *Hepatology* 2019; 69: 2672-2682.
- [3] Lonardo A, Nascimbeni F, Maurantonio M, Marrazzo A, Rinaldi L and Adinolfi LE. Nonalcoholic fatty liver disease: evolving paradigms. *World J Gastroenterol* 2017; 23: 6571-6592.
- [4] Fabbrini E, Sullivan S and Klein S. Obesity and nonalcoholic fatty liver disease: biochemical, metabolic, and clinical implications. *Hepatology* 2010; 51: 679-689.
- [5] Younossi Z, Anstee QM, Marietti M, Hardy T, Henry L, Eslam M, George J and Bugianesi E. Global burden of NAFLD and NASH: trends, predictions, risk factors and prevention. *Nat Rev Gastroenterol Hepatol* 2018; 15: 11-20.
- [6] Estes C, Anstee QM, Arias-Loste MT, Bantel H, Bellentani S, Caballeria J, Colombo M, Craxi A, Crespo J, Day CP, Eguchi Y, Geier A, Kondili LA, Kroy DC, Lazarus JV, Loomba R, Manns MP,

A predictive model for NAFLD

- Marchesini G, Nakajima A, Negro F, Petta S, Ratziu V, Romero-Gomez M, Sanyal A, Schattenberg JM, Tacke F, Tanaka J, Trautwein C, Wei L, Zeuzem S and Razavi H. Modeling NAFLD disease burden in China, France, Germany, Italy, Japan, Spain, United Kingdom, and United States for the period 2016-2030. *J Hepatol* 2018; 69: 896-904.
- [7] Chalasani N, Younossi Z, Lavine JE, Diehl AM, Brunt EM, Cusi K, Charlton M and Sanyal AJ. The diagnosis and management of non-alcoholic fatty liver disease: practice guideline by the American Association for the Study of Liver Diseases, American College of Gastroenterology, and the American Gastroenterological Association. *Hepatology* 2012; 55: 2005-2023.
- [8] Ma X, Liu S, Zhang J, Dong M, Wang Y, Wang M and Xin Y. Proportion of NAFLD patients with normal ALT value in overall NAFLD patients: a systematic review and meta-analysis. *BMC Gastroenterol* 2020; 20: 10.
- [9] Poynard T, Ratziu V, Naveau S, Thabut D, Charlotte F, Messous D, Capron D, Abella A, Massard J, Ngo Y, Munteanu M, Mercadier A, Manns M and Albrecht J. The diagnostic value of biomarkers (SteatoTest) for the prediction of liver steatosis. *Comp Hepatol* 2005; 4: 10.
- [10] Li W, Huang Y, Zhuang BW, Liu GJ, Hu HT, Li X, Liang JY, Wang Z, Huang XW, Zhang CQ, Ruan SM, Xie XY, Kuang M, Lu MD, Chen LD and Wang W. Multiparametric ultrasonomics of significant liver fibrosis: a machine learning-based analysis. *Eur Radiol* 2019; 29: 1496-1506.
- [11] Liu X, Hou Y, Wang X, Yu L, Wang X, Jiang L and Yang Z. Machine learning-based development and validation of a scoring system for progression-free survival in liver cancer. *Hepatol Int* 2020; 14: 567-576.
- [12] Chalasani N, Younossi Z, Lavine JE, Charlton M, Cusi K, Rinella M, Harrison SA, Brunt EM and Sanyal AJ. The diagnosis and management of nonalcoholic fatty liver disease: practice guidance from the American Association for the Study of Liver Diseases. *Hepatology* 2018; 67: 328-357.
- [13] Abeel T, Helleputte T, Van de Peer Y, Dupont P and Saeys Y. Robust biomarker identification for cancer diagnosis with ensemble feature selection methods. *Bioinformatics* 2010; 26: 392-398.
- [14] Cai Z, Xu D, Zhang Q, Zhang J, Ngai SM and Shao J. Classification of lung cancer using ensemble-based feature selection and machine learning methods. *Mol Biosyst* 2015; 11: 791-800.
- [15] He Z and Yu W. Stable feature selection for biomarker discovery. *Comput Biol Chem* 2010; 34: 215-225.
- [16] Hilbe JM. Logistic regression models. London: CRC press; 2009.
- [17] Breiman L. Random forests. In: Schapire, editor. *Machine Learning*. The Netherlands; Kluwer Academic Publishers: 2002. pp. 5-32.
- [18] Chang CC and Lin CJ. LIBSVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* 2011; 2: 1-39.
- [19] Perrier E, Rondeau P, Poupin M, Le Bellego L, Armstrong LE, Lang F, Stookey J, Tack I, Vergne S and Klein A. Relation between urinary hydration biomarkers and total fluid intake in healthy adults. *Eur J Clin Nutr* 2013; 67: 939-943.
- [20] Verslype C. Evaluation of abnormal liver-enzyme results in asymptomatic patients. *Acta Clin Belg* 2004; 59: 285-289.
- [21] Nanji AA, French SW and Freeman JB. Serum alanine aminotransferase to aspartate aminotransferase ratio and degree of fatty liver in morbidly obese patients. *Enzyme* 1986; 36: 266-269.
- [22] Long MT, Pedley A, Colantonio LD, Massaro JM, Hoffmann U, Muntner P and Fox CS. Development and validation of the framingham steatosis index to identify persons with hepatic steatosis. *Clin Gastroenterol Hepatol* 2016; 14: 1172-1180, e1172.
- [23] Tarantino G and Finelli C. What about non-alcoholic fatty liver disease as a new criterion to define metabolic syndrome? *World J Gastroenterol* 2013; 19: 3375-3384.
- [24] Chao TT, Hsieh CH, Lin JD, Wu CZ, Hsu CH, Pei D, Chen YL, Liang YJ and Chang JB. Use of white blood cell counts to predict metabolic syndrome in the elderly: a 4 year longitudinal study. *Aging Male* 2014; 17: 230-237.
- [25] Yang H, Fu YQ, Yang B, Zheng JS, Zeng XY, Zeng W, Fan ZF, Chen M, Wang L and Li D. Positive association between the metabolic syndrome and white blood cell counts in Chinese. *Asia Pac J Clin Nutr* 2017; 26: 141-147.
- [26] Riley LK and Rupert J. Evaluation of patients with leukocytosis. *Am Fam Physician* 2015; 92: 1004-1011.
- [27] Wang S, Zhang C, Zhang G, Yuan Z, Liu Y, Ding L, Sun X, Jia H and Xue F. Association between white blood cell count and non-alcoholic fatty liver disease in urban Han Chinese: a prospective cohort study. *BMJ Open* 2016; 6: e010342.
- [28] Lee YJ, Lee HR, Shim JY, Moon BS, Lee JH and Kim JK. Relationship between white blood cell count and nonalcoholic fatty liver disease. *Dig Liver Dis* 2010; 42: 888-894.
- [29] Chung GE, Yim JY, Kim D, Kwak MS, Yang JI, Chung SJ, Yang SY and Kim JS. Associations between white blood cell count and the development of incidental nonalcoholic fatty liver

A predictive model for NAFLD

- disease. *Gastroenterol Res Pract* 2016; 2016: 7653689.
- [30] Stefaniuk P, Cianciara J and Wiercinska-Drapalo A. Present and future possibilities for early diagnosis of hepatocellular carcinoma. *World J Gastroenterol* 2010; 16: 418-424.
- [31] Tangkijvanich P, Tosukhowong P, Bunyongyod P, Lertmaharit S, Hanvivatvong O, Kullavanijaya P and Poovorawan Y. Alpha-L-fucosidase as a serum marker of hepatocellular carcinoma in Thailand. *Southeast Asian J Trop Med Public Health* 1999; 30: 110-114.
- [32] Ramm GA and Ruddell RG. Hepatotoxicity of iron overload: mechanisms of iron-induced hepatic fibrogenesis. *Semin Liver Dis* 2005; 25: 433-449.
- [33] Yajima D, Motani H, Hayakawa M, Sato Y, Sato K and Iwase H. The relationship between cell membrane damage and lipid peroxidation under the condition of hypoxia-reoxygenation: analysis of the mechanism using antioxidants and electron transport inhibitors. *Cell Biochem Funct* 2009; 27: 338-343.
- [34] Lu ZY, Cen C, Shao Z, Chen XH, Xu CF and Li YM. Association between serum alpha-L-fucosidase and non-alcoholic fatty liver disease: cross-sectional study. *World J Gastroenterol* 2016; 22: 1884-1890.
- [35] Dongiovanni P, Valenti L, Ludovica Fracanzani A, Gatti S, Cairo G and Fargion S. Iron depletion by deferoxamine up-regulates glucose uptake and insulin signaling in hepatoma cells and in rat liver. *Am J Pathol* 2008; 172: 738-747.
- [36] Bai CH, Wu MS, Owaga E, Cheng SY, Pan WH and Chang JS. Relationship between hemoglobin levels and risk for suspected non-alcoholic fatty liver in Taiwanese adults. *Chin J Physiol* 2014; 57: 286-294.
- [37] Chung GE, Yim JY, Kim D, Kwak MS, Yang JI, Chung SJ, Yang SY and Kim JS. Associations between hemoglobin concentrations and the development of incidental metabolic syndrome or nonalcoholic fatty liver disease. *Dig Liver Dis* 2017; 49: 57-62.
- [38] Giorgio V, Mosca A, Alterio A, Alisi A, Grieco A, Nobili V and Miele L. Elevated hemoglobin level is associated with advanced fibrosis in pediatric nonalcoholic fatty liver disease. *J Pediatr Gastroenterol Nutr* 2017; 65: 150-155.
- [39] Juárez-Hernández E, C Chávez-Tapia N, C Brihueza-Alcántara D, Uribe M, H Ramos-Ostos M and Nuño-Lámbarri N. Association between serum hemoglobin levels and non alcoholic fatty liver disease in a mexican population. *Ann Hepatol* 2018; 17: 577-584.
- [40] Whitfield JB. Gamma glutamyl transferase. *Crit Rev Clin Lab Sci* 2001; 38: 263-355.
- [41] Satapati S, Kucejova B, Duarte JA, Fletcher JA, Reynolds L, Sunny NE, He T, Nair LA, Livingston KA, Fu X, Merritt ME, Sherry AD, Malloy CR, Shelton JM, Lambert J, Parks EJ, Corbin I, Magnuson MA, Browning JD and Burgess SC. Mitochondrial metabolism mediates oxidative stress and inflammation in fatty liver. *J Clin Invest* 2015; 125: 4447-4462.
- [42] Jarcuska P, Janicko M, Drazilova S, Senajova G, Veseliny E, Fedacko J, Siegfried L, Kristian P, Tkac M Jr, Pella D, Marekova M, Geckova AM and Jarcuska P; HepaMeta Team. Gamma-glutamyl transpeptidase level associated with metabolic syndrome and proinflammatory parameters in the young Roma population in eastern Slovakia: a population-based study. *Cent Eur J Public Health* 2014; 22 Suppl: S43-50.
- [43] Li Q, Lu C, Li W, Huang Y and Chen L. The gamma-glutamyl transpeptidase to platelet ratio for non-invasive assessment of liver fibrosis in patients with chronic hepatitis B and non-alcoholic fatty liver disease. *Oncotarget* 2017; 8: 28641-28649.
- [44] Svegliati-Baroni G, Pierantonelli I, Torquato P, Marinelli R, Ferreri C, Chatgialiloglu C, Bartolini D and Galli F. Lipidomic biomarkers and mechanisms of lipotoxicity in non-alcoholic fatty liver disease. *Free Radic Biol Med* 2019; 144: 293-309.
- [45] Kimura T, Tanaka N, Fujimori N, Yamazaki T, Katsuyama T, Iwashita Y, Pham J, Joshita S, Pydi SP and Umemura T. Serum thrombospondin 2 is a novel predictor for the severity in the patients with NAFLD. *Liver Int* 2021; 41: 505-514.
- [46] Mele C, Crinò A, Fintini D, Mai S, Convertino A, Bocchini S, Di Paolo P, Grugni G, Aimaretti G, Scacchi M and Marzullo P. Angiotensin-like 8 (ANGPTL8) as a potential predictor of NAFLD in paediatric patients with Prader-Willi Syndrome. *J Endocrinol Invest* 2021; 44: 1447-1456.

A predictive model for NAFLD

Supplementary Table 1. All the clinical characteristics of subjects with and without NAFLD in the training cohort

Characters	NAFLD N=53	No NAFLD N=45	P value
Gender (male/female, n)	32/21	25/20	0.02
ALT (U/L)	31.21±21.54	19.30±12.93	<0.001
AST (U/L)	22.24±4.61	20.73±12.83	0.27
ALT to AST	1.33±0.64	0.90±0.29	<0.001
GGT (U/L)	35.88±29.93	21.95±13.96	0.02
AFU (U/L)	26.45±7.91	23.32±7.12	0.046
WBC (× 10 ⁹ /L)	6.40±1.70	5.30±1.31	<0.001
TG (mmol/L)	1.83±1.26	1.25±0.53	<0.001
Hb (g/L)	144.98±11.69	138.22±15.76	0.03
NEUT (10 ⁹ /L)	3.40±1.11	2.90±0.94	0.02
NEUT%	52.55±8.96	53.13±7.82	0.74
LYM (10 ⁹ /L)	2.78±2.52	2.00±0.62	0.03
LYM%	37.13±9.28	37.17±8.60	0.98
MON (10 ⁹ /L)	0.40±0.12	0.34±0.12	0.02
MON%	6.37±1.62	6.46±1.65	0.78
EOS (10 ⁹ /L)	0.16±0.10	0.15±0.11	0.78
EOS%	2.49±1.59	2.77±1.59	0.39
BASO (10 ⁹ /L)	0.04±0.02	0.03±0.02	0.04
BASO%	0.57±0.26	0.55±0.30	0.75
RBC (10 ¹² /L)	4.63±0.46	4.47±0.43	0.09
HCT (%)	43.12±3.60	41.56±4.08	0.05
MCV (fL)	92.55±4.07	93.16±5.07	0.52
MCH (pg)	31.07±1.53	31.09±2.00	0.96
MCHC (g/L)	333.40±17.47	333.73±11.43	0.91
RDW (fL)	42.71±5.39	44.18±3.56	0.12
RDW (%)	12.56±0.53	12.79±1.23	0.24
PLT (fL)	229.98±51.26	224.73±50.05	0.62
PCT (%)	0.22±0.05	0.23±0.09	0.74
MPV (fL)	9.58±1.13	9.62±1.13	0.89
P-LCR%	27.46±6.16	27.81±7.04	0.80
PDW (%)	13.99±2.13	14.10±2.37	0.81
TP (g/L)	71.75±6.31	71.14±10.90	0.74
ALB (g/L)	43.50±3.50	43.43±3.42	0.93
GLOB (g/L)	28.43±4.77	29.27±4.88	0.41
A/G	1.58±0.29	1.54±0.27	0.48
TB (μmol/L)	13.31±3.85	13.20±4.36	0.90
DB (μmol/L)	2.23±0.78	2.16±0.69	0.67
IB (μmol/L)	11.08±3.25	11.10±3.86	0.97
ALP (U/L)	78.94±19.56	73.03±21.74	0.18
TBA (μmol/L)	4.89±3.08	5.23±4.17	0.66
GLU (mmol/L)	5.00±0.80	4.80±0.86	0.27
BUN (mmol/L)	5.45±1.53	5.10±1.02	0.20
CRE (μmol/L)	68.38±13.98	74.01±49.60	0.47
BUN/CRE	0.08±0.02	0.08±0.02	0.82
UA (μmol/L)	360.34±86.08	312.26±70.99	0.00
CK (U/L)	118.83±116.29	96.49±83.32	0.31

A predictive model for NAFLD

CK-MB (U/L)	15.30±12.95	14.63±7.09	0.77
LDH (mmol/L)	178.06±27.60	170.05±34.31	0.24
HBDH (U/L)	121.39±18.69	119.26±25.10	0.66
TC (mmol/L)	5.32±1.01	5.11±1.08	0.35
HDL (mmol/L)	1.24±0.28	1.39±0.31	0.02
LDL (mmol/L)	3.19±0.72	3.05±0.80	0.39
APO A1 (g/L)	1.53±0.52	1.45±0.30	0.41
APO B (g/L)	1.01±0.46	0.98±0.26	0.26
APO A1/APO B	1.35±0.39	1.62±0.70	0.04
LP (a) (mg/dL)	18.80±17.30	23.60±22.82	0.29

Supplementary Table 2. All the clinical characteristics of subjects with and without NAFLD in the validation cohort

Characters	NAFLD group N=81	Control group N=87	P value
ALT (U/L)	31.45±19.24	23.27±15.86	<0.001
AST (U/L)	21.82±8.52	19.22±6.69	0.03
ALT to AST	1.41±0.44	1.16±0.52	<0.001
GGT (U/L)	33.78±14.53	18.20±6.47	0.02
AFU (U/L)	29.35±7.88	26.31±7.05	0.022
WBC (× 10 ⁹ /L)	6.46±1.69	5.80±1.25	0.012
TG (mmol/L)	2.03±1.89	1.39±1.14	0.022
Hb (g/L)	143.54±14.10	138.80±17.20	0.08
LYM (10 ⁹ /L)	2.19±0.75	1.91±0.59	0.01
MON (10 ⁹ /L)	0.46±0.15	0.37±0.13	<0.001
MON%	7.32±2.02	6.42±2.02	<0.001
BASO%	0.53±0.25	0.44±0.28	0.04
RDW (%)	12.37±0.55	22.65±13.79	<0.001
MPV (fL)	9.88±0.94	10.06±0.95	0.23
P-LCR%	24.42±6.47	25.82±7.39	0.20
PDW (%)	12.78±2.31	13.14±2.58	0.30
TP (g/L)	69.34±5.47	69.85±5.15	0.55
ALB (g/L)	43.14±3.83	42.96±3.59	0.77
GLOB (g/L)	26.29±4.26	26.83±3.65	0.39
A/G	1.68±0.31	1.64±0.28	0.30
TB (μmol/L)	15.80±7.74	14.20±5.96	0.14
DB (μmol/L)	4.27±2.04	4.07±1.84	0.51
IB (μmol/L)	11.40±6.02	10.12±4.55	0.13
ALP (U/L)	63.82±17.39	62.49±27.33	0.71
TBA (μmol/L)	4.09±3.09	4.22±3.71	0.82
GLU (mmol/L)	5.49±1.27	5.00±1.27	0.40
BUN (mmol/L)	5.35±1.20	5.13±1.24	0.40
CRE (μmol/L)	76.34±18.13	77.37±16.35	0.71
BUN/CRE	14.58±8.35	14.61±7.52	0.98
UA (μmol/L)	366.70±78.65	306.36±80.34	0.00
CK (U/L)	113.28±83.31	90.27±35.60	0.05
CK-MB (U/L)	12.19±4.57	13.25±5.26	0.23
LDH (mmol/L)	162.86±33.33	159.17±37.33	0.53
HBDH (U/L)	128.30±22.32	124.18±27.89	0.59
TC (mmol/L)	4.54±0.54	4.15±1.60	<0.001

A predictive model for NAFLD

HDL (mmol/L)	1.33±0.30	1.46±0.36	0.04
LDL (mmol/L)	3.21±0.75	2.92±0.87	0.07
APO A1 (g/L)	1.38±0.24	1.52±0.22	0.01
APO B (g/L)	1.00±0.24	0.85±0.19	0.01
APO A1/APO B	1.48±0.51	1.88±0.53	<0.001
LP (a) (mg/dL)	18.55±16.30	20.64±21.22	0.37

Supplementary Table 3. Logical regression coefficients of the candidate Biomarkers in the training cohort

Name	Coefficients
(Intercept)	-8.77
WBC	0.29
ALT to AST	2.23
AFU	0.07
Hb	0.01
TG	0.72
GGT	0.01

Supplementary Table 4. The best cutoff value for the predicted model in the training cohort

	Threshold	AUROC	Specificity	Sensitivity
Best	0.53	0.821	0.733	0.765

Supplementary Table 5. Characteristics of mice fed with HFD or CD

Characters	HFD group N=32	CD group N=28	P value
ALT (U/L)	48.14±15.57	45.48±11.02	0.58
AST (U/L)	193.77±63.35	182.32±50.56	0.59
ALT to AST	0.27±0.10	0.26±0.06	0.78
GGT (IU/L)	6.29±0.90	6.14±0.60	0.58
AFU (U/L)	6.14±0.71	6.02±0.47	0.58
WBC (× 10 ⁹ /L)	9.38±1.08	8.98±1.05	0.33
TG (mmol/L)	0.90±0.12	0.85±0.11	0.24
Hb (g/L)	140.07±15.99	129.03±20.32	0.11