

Chinese forest fire occurrence prediction based on machine learning methods

yudong Li

Beijing Forestry University

Zhongke Feng (✉ zhongkefeng@bjfu.edu.cn)

Beijing Forestry University <https://orcid.org/0000-0003-1602-5045>

Ziyu Zhao

Beijing Forestry University

Wenyuan Ma

Beijing Forestry University

Shilin Chen

Beijing Forestry University

Hanyue Zhang

Beijing Forestry University

Research

Keywords: forest fire occurrence in China, feature selection, forest fire driving factors, machine learning, prediction model, forest fire prevention and control

Posted Date: January 12th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-62305/v2>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 [Title Page]

2 **Chinese forest fire occurrence prediction based on machine learning methods**

3 **Yudong Li,¹ Zhongke Feng,¹ Ziyu Zhao,¹Wenyuan Ma,¹Shilin Chen,¹Hanyue Zhang,¹**

4 **1 Precision Forestry Key Laboratory of Beijing, Beijing Forestry University, Beijing 100083,**

5 **China.**

6 **E-mail address :**

7 **Yudong Li:** lyd85842@163.com; **Zhongke Feng:** zhongkefeng@bjfu.edu.cn

8 **Ziyu Zhao:** zhaozy0315@126.com; **Wenyuan Ma:**imawenyuan@126.com;

9 **Shilin Chen:** chenshilin@bjfu.edu.cn; **Hanyue Zhang:** hanyue.zhang@foxmail.com

10 **Correspondence should be addressed to Zhongke Feng:** zhongkefeng@bjfu.edu.cn

11 **Permanent address: Beijing Forestry University**

12 [Manuscript]

13 **Chinese forest fire occurrence prediction based on machine**
14 **learning methods**

15 Yudong Li,¹ Zhongke Feng,¹ Ziyu Zhao,¹ Wenyuan Ma,¹Shilin Chen,¹Hanyue Zhang,¹

16 1 Precision Forestry Key Laboratory of Beijing, Beijing Forestry University, Beijing 100083, China.

17 Correspondence should be addressed to Zhongke Feng: zhongkefeng@bjfu.edu.cn

18 **Abstract**

19 Forest fires can cause serious harm. Scientifically predicting forest fires is an important basis for

20 preventing them. Currently, there is little research on the prediction of long time-series forest fires in

21 China. Choosing a suitable forest fire prediction model and predicting the probability of Chinese forest
22 fire occurrence are of great importance to China's forest fire prevention and control work. Based on fire
23 hotspot, meteorological, terrain, vegetation, infrastructure, and socioeconomic data collected from 2003
24 to 2016, we used a random forest model as a feature-selection method to identify 13 major drivers of
25 forest fires in China. The forest fire prediction models developed in this study are based on four machine-
26 learning algorithms: an artificial neural network, a radial basis function network, a support-vector
27 machine, and a random forest. The models were evaluated using the five performance indicators of
28 accuracy, precision, recall, f1 value, and area under the curve. We used the optimal model to obtain the
29 probability of forest fire occurrence in various provinces in China and created a spatial distribution map
30 of the areas with high incidences of forest fires. The results showed that the prediction accuracy of the
31 four forest fire prediction models was between 75.8% and 89.2%, and the area under the curve value was
32 between 0.840 and 0.960. The random forest model had the highest accuracy (89.2%) and area under the
33 curve value (0.96); thus, it was used as the optimal model to predict the probability of forest fire
34 occurrence in China. The prediction results indicate that the areas with high incidences of forest fires are
35 mainly concentrated in north-eastern China (Heilongjiang Province and northern Inner Mongolia
36 Autonomous Region) and south-eastern China (including Fujian Province and Jiangxi Province). In areas
37 at high risk of forest fire, management departments can improve forest fire prevention and control by
38 establishing watch towers and using other monitoring equipment. This study helps in understanding the
39 main drivers of forest fires in China, provides a reference for the selection of high-precision forest fire
40 prediction models, and provides a scientific basis for China's forest fire prevention and control work.

41 **Keywords:** forest fire occurrence in China; feature selection; forest fire driving factors; machine learning;
42 prediction model; forest fire prevention and control

43 **1. Introduction**

44 Forest fires are one of the most dangerous natural disasters. They have received worldwide attention due
45 to their rapid spread, their low controllability, and the hazards they pose[1][2]. Forest fires have varying
46 degrees of impact on human health and safety, the ecological environment and resources, and society and
47 the economy[3]. Forest fire prevention has therefore become a key research topic in the fields of forestry
48 and ecology[4][5][6].

49 The most effective way to control forest fires is to detect them quickly. Detection is usually divided into
50 three categories: satellite monitoring, smoke detection, and local perception (such as data analysis).
51 Satellite monitoring is expensive, is affected by delays, and is not fully applicable to all locations[7].
52 Smoke detection also requires expensive equipment and maintenance work. In contrast, data for forest
53 fire analysis are easy to obtain, and data analysis is less expensive than the first two methods[8].

54 In recent years, researchers have mostly used data analysis methods to obtain the degree of risk of forest
55 fires in the study area[9][10][11]. Many researchers have used geospatial technologies such as GIS and
56 remote sensing (RS) and have used data analysis models to analyse and evaluate an area's sensitivity to
57 forest fire[12][13][14]. Forest fire susceptibility is the probability of wildfire occurrence based on certain
58 thresholds[15]. When the probability of wildfire is high, the potential forest fire risk is high. By
59 establishing a forest fire prediction model, we can predict the probability of the occurrence of a forest
60 fire and then strictly manage the area where the fire is likely to occur. This approach can directly reduce
61 the occurrence of forest fires as well as the associated potential casualties and economic losses[16][17].
62 This method is therefore of great significance in forest fire prediction and prevention[18]

63 Much research has been conducted on forest fire prediction models. Logistic regression models are the

64 most commonly used models, and they have the advantage of solving the classification problem[19][20].
65 In recent years, geographically weighted regression models have also been used[21][22]. This method
66 can provide a reasonable explanation for spatial heterogeneity, but the regression analyses can be
67 performed only on continuous variables; the method lacks analysis of categorical variables. Similarly,
68 many researchers have used generalized linear regression models for forest fire prediction. Liao et al.
69 (2008) used the zero-inflated Poisson model to predict the frequency of forest fires in Japan in 2000 [23].
70 Mandallaz et al. (1997) used the Poisson model to predict forest fires in France, Italy, etc. [24]. Guo et
71 al. (2010) used ordinary least squares regression, zero-inflated negative binomial regression and the zero-
72 inflated negative binomial model to predict the number of forest fires in the Greater Xing'an Mountains
73 area of Heilongjiang Province, China[25].

74 The development of artificial intelligence has led researchers to focus on building a forest fire prediction
75 model using machine learning algorithms[26][27][28][29][30][31][32][33][34][35]. Researchers have
76 used machine learning models to predict wildfires with greater accuracy and faster automation. For
77 example, Camp et al.(1997)used decision trees (DTs) to identify historic forest fire shelters[36];
78 Pourghasemi et al. (2016)used the Mamdani fuzzy logic model to evaluate forest fire prediction
79 ability[37]; and Li et al. (2020)used the long-term memory (LSTM) neural network model to study the
80 burning area of wildfires[38]. For forest fire prediction, researchers have used more machine learning
81 methods, such as artificial neural network models and support vector machine models.

82 Artificial neural networks (ANNs) consist of neurons with adjustable connection weights. Compared
83 with traditional multiple linear regression models or parametric regression models, neural networks have
84 better self-organization and self-learning capabilities, and they have been widely used in forest fire

85 prediction[39][40][41]. For example, Maeda et al. (2009) used ANNs and multitemporal images from
86 MODIS/Terra-Aqua sensors to detect areas at high risk of forest fires in the Amazon region of Brazil
87 [42]. The results showed that the error was small, and the predictions were accurate. Sakr et al. (2011)
88 predicted the occurrence of forest fires in developing countries through two meteorological factors using
89 artificial neural networks [43]. A radial basis function (RBF) neural network is a three-layer neural
90 network, and it is a special case of a back-propagation neural network. At present, little research has used
91 RBF neural networks for forest fire prediction. Samaher (2018) used an RBF neural network to predict
92 the forest fire risk in natural parks in Portugal [27]. Support-vector machines (SVMs) are most suitable
93 for the binary classification of data in the form of supervised learning. SVMs apply the principle of
94 structural risk minimization and have good learning ability. In recent years, researchers have begun to
95 use SVMs to predict forest fires[8][44][45][46]. Samaher (2018) used five different soft computing (SC)
96 technologies, including an SVM algorithm, to predict areas at risk of forest fires[27]. He determined that
97 the SVM algorithm provides more accurate predictions than the other four algorithms. Cortez et al. (2007)
98 used five different data mining (DM) algorithms to predict the area at risk of forest fires in the north-
99 eastern region of Portugal[8]. Their results showed that the prediction effect of the model was good.
100 Based on Cortez's research, Xu et al. (2012) used the semidefinite programming model to select the
101 optimal kernel function of the SVM to establish an SVM model for forest fire prediction [7]. The mean
102 square error was small, and the model effect was good. The random forest (RF) algorithm is a well-
103 known integrated learning algorithm that can provide higher accuracy than other algorithms. At present,
104 the use of RFs to predict forest fires is relatively established[47][48][49]. Liang et al. (2016) used an RF
105 model to predict the occurrence of forest fires in Fujian Province, China, with an accuracy rate of 85%
106 [50]. Pourtaghi et al. (2016) used an RF algorithm to study the sensitivity of forest fires in Golestan

107 Province, Iran, and their results showed that the model achieved the desired accuracy[51].
108 Most of the current research focuses on certain areas, and there are few studies on the prediction and
109 analysis of long time-series in China. Most studies have concentrated on the temporal and spatial changes
110 and influencing factors of forest fires in specific years[52][53][54][55][56]. The results of previous
111 research are therefore localized and limited, and there is a lack of research investigating the most suitable
112 and high-precision forest fire prediction model on the national scale.

113 In this study, we selected a variety of forest fire driving factors to build four prediction models based on
114 machine learning algorithms. The models were evaluated using data on Chinese forest fires from 2003
115 to 2016. The study has three objectives: (1) identify the main forest fire driving factors and their impacts
116 in China; (2) select the most suitable model for forest fire prediction in China by creating four models
117 and comparing and analysing the fitting results; and (3) use the model that offers the most accurate
118 predictions to create a forest fire probability map for China and put forward recommendations for forest
119 fire prevention.

120 **2. Materials and Methods**

121 **2.1 Study Area and Data Resources**

122 Located in East Asia on the west coast of the Pacific Ocean, China's territory is vast, with a total land
123 area of approximately 9.6 million square kilometres. The topography is high in the west, with vast
124 mountains and plateaus, and low in the east. The distance between the east and the west of the country is
125 approximately 5,000 kilometres; the coastline of the mainland is more than 18,000 kilometres in length;
126 and the temperature and precipitation are diverse, forming a variety of climates. The distribution of forest

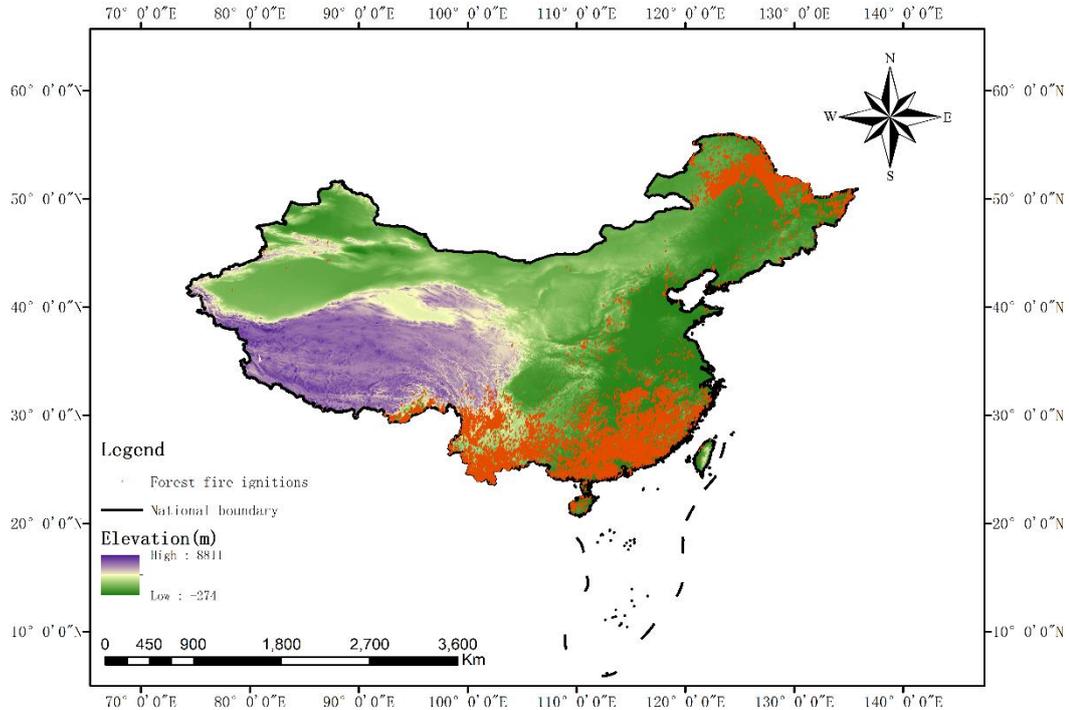
127 resources in China is uneven and is mainly distributed in the north-eastern, southern, and south-western
128 regions. The forested area is 220 million hectares, and the forest coverage rate is 22.96%.

129 The research data were divided into six parts: fire ignition data, meteorological data, terrain data,
130 vegetation data, infrastructure data, and socioeconomic data[70]. The fire point data were derived from
131 NASA's *Global Fire Atlas with Characteristics of Individual Fires, 2003–2016* (<https://daac.ornl.gov/>).

132 The *Global Fire Atlas* is a global dataset that tracks the daily dynamics of single fires. For each individual
133 fire, the dataset provides information about the fire's timing and location, scale, perimeter, duration,
134 speed, and direction of spread. These individual fire characteristics are based on the *Global Fire Atlas*
135 algorithms and estimated combustion day information from a 500-metre resolution product of the 6
136 MCD64A1 combustion zone product of the Medium Resolution Imaging Spectroradiometer (MODIS)
137 collection.

138 This study used fire point data for forest land in China from 2003 to 2016. The final number of fire point
139 data points was 32,746 (excluding Taiwan). The meteorological data were derived from the 14-day daily
140 value dataset of the China Meteorological Data Network (<http://data.cma.cn/dataService/>). The dataset
141 includes eight elements, such as the barometric pressure, temperature, relative humidity, and precipitation
142 at the station. Digital elevation model (DEM) data were obtained through the Geospatial Data Cloud
143 website (<http://www.gscloud.cn/>). Vegetation data were represented by the normalized difference
144 vegetation index (NDVI), and the spatial distribution dataset of China's Quarterly Vegetation Index came
145 from the Resource and Environment Data Cloud Platform (<http://www.resdc.cn/>). The basic geographic
146 data were taken from the "National Basic Geographic Database of 1:1 Million" on the website of the
147 National Geographic Information Resources Directory System. The data include the locations of railways,

148 highways, water systems, and residential areas. The socioeconomic data include population density and
149 GDP per capita, and the grid data of the spatial distribution of population and GDP were obtained from
150 the Resource and Environment Data Cloud Platform. Figure 1 shows the map of the study area.



151

152

Fig. 1 Map of the study area

153 **2.2 Data Pre-processing**

154 **2.2.1 Variable Handling**

155 The dependent variable is a binary variable (i.e., whether or not a forest fire occurs); thus, we used
156 ArcGIS 10.4 to create a certain percentage of random points (non-fire points) and assigned a value of 1
157 to fire points and a value of 0 to non-fire points[57]. To ensure that the data were not overdispersed,
158 random points were selected according to experience in a ratio of 1:1[58], and in principle, randomness
159 in space and time should be followed[59]. We used ArcGIS 10.4 software to create random points and

160 then used the 2015 national land-use data as a basis to exclude random points that were located in bodies
161 of water or urban land. We obtained a total of 65,492 fire points and random points.

162 For the meteorological data, we first used ArcGIS 10.4 to match the sample points with the nearest
163 meteorological station using the Thiessen polygon method. We then extracted the corresponding sample
164 point weather data and used an SQL server database to match the daily weather data. For the terrain data,
165 we used the spatial analysis tool in ArcGIS 10.4 to extract the slope and aspect of the obtained DEM data.
166 Seasonal climatic differences have an impact on vegetation status; thus, we divided the year into spring
167 (March, April, May), summer (June, July, August), autumn (September, October, November), and winter
168 (December, January, February)[60]. We used the extraction and analysis tools of ArcGIS to extract the
169 NDVI data for the sample points on an annual basis and a quarterly basis.

170 Similarly, from the infrastructure data and socioeconomic data, we extracted the information
171 corresponding to the sample points. We set the aspect and special festivals as categorical variables and
172 the others as continuous variables. Table 1 shows the classification of aspect. During certain traditional
173 festivals in China, people burn paper to commemorate their loved ones, increasing the probability of a
174 forest fire. We classified (value 1) the following dates of these events as special festivals: Chinese New
175 Year's Eve, the first day of the first lunar month, the second day of the first lunar month, the fifteenth
176 day of the first lunar month, and Qingming Festival and Zhongyuan Festival (July 15th of the lunar
177 calendar). Non-special festivals were set to 0.

178 **Table 1: Descriptions of aspect classifications**

Aspect	Azimuth (degree)	Classification
--------	------------------	----------------

Gentle Slope	-1	0
Shady Slope	0~67.5, 337.5~360	1
Semi-shady Slope	67.5~112.5, 292.5~337.5	2
Sunny Slope	157.5~247.5	3
Semi-sunny Slope	112.5~157.5, 247.5~292.5	4

179 After processing, we obtained 20 independent variables and their possible values (see Table 2). Finally,
180 we performed data cleaning on the sample points and the various types of data extracted to remove
181 abnormal samples from the original dataset (including some samples with missing data and samples with
182 observations that were significantly outside the normal range).

183

Table 2: Descriptions of independent variables

Category	Independent Variable	Symbol	Variable Type
Location	Longitude (°)	<i>Lon</i>	Continuous Variable
	Latitude (°)	<i>Lat</i>	Continuous Variable
Terrain	Altitude (m)	<i>Alt</i>	Continuous Variable
	Slope (°)	<i>Slo</i>	Continuous Variable
	Aspect	<i>Asp</i>	Categorical Variable
Meteorology	Average Surface Temperature (°C)	<i>Avst</i>	Continuous Variable
	Daily Maximum Surface temperature (°C)	<i>Mast</i>	Continuous Variable
	Cumulative Precipitation at 20–20 (mm)	<i>Pre</i>	Continuous Variable
	Average Relative Humidity (%)	<i>Arh</i>	Continuous Variable
	Hours of Sunshine (h)	<i>Suh</i>	Continuous Variable

	Average Temperature (°C)	<i>Ate</i>	Continuous Variable
	Daily Maximum Temperature (°C)	<i>Mate</i>	Continuous Variable
	Average Wind Speed (m/s)	<i>Aws</i>	Continuous Variable
	Maximum Wind Speed (m/s)	<i>Mws</i>	Continuous Variable
Infrastructure	Distance from Fire Point to Highway (m)	<i>Hig</i>	Continuous Variable
	Closest Distance from Fire Point to Residential Area (m)	<i>Set</i>	Continuous Variable
Social humanity	Population	<i>Pop</i>	Continuous Variable
	GDP	<i>GDP</i>	Continuous Variable
	Special Festival	<i>Sfe</i>	Categorical Variable
Vegetation	NDVI	<i>NDVI</i>	Continuous Variable

184 2.2.2 Data Normalization

185 Given the different dimensions and magnitudes of the factors above, the data were normalized to
186 eliminate the variation in dimensions, avoid large differences in the magnitudes of the input and output
187 data, and balance the contributions of various factors. All data were converted to values between 0 and
188 1. Table 3 shows the normalized formulas and specific interpretations of the independent variables.

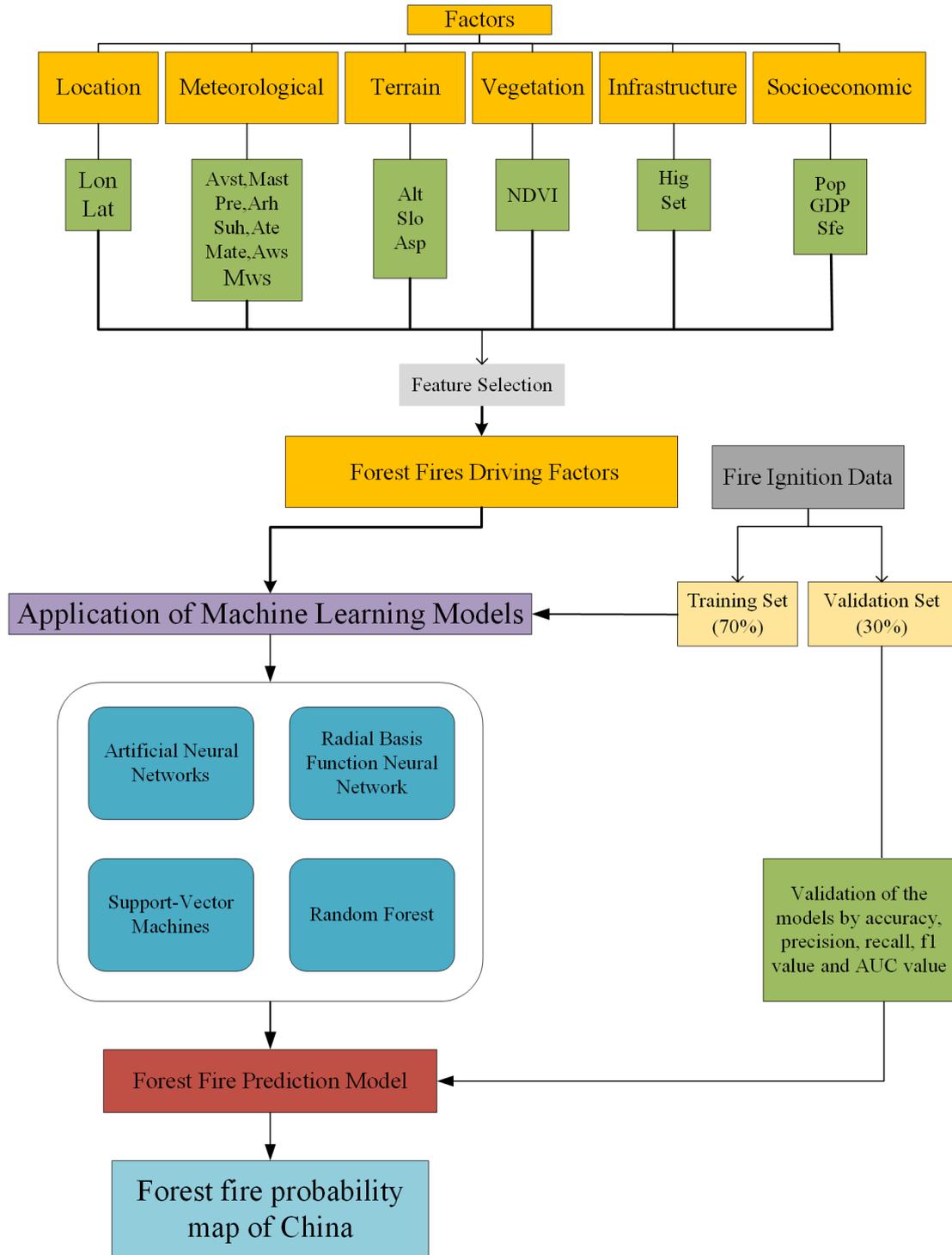
189 **Table 3: Normalized formulas and explanations**

No.	Formula	Explanation	Variables using this formula
(1)	$x_i^* = \frac{x_i - x_{min}}{x_{max} - x_{min}}$	x_i and x_i^* are the values before and after data normalization, respectively; x_{max} and x_{min} are the maximum and minimum values of the full sample data, respectively.	<i>Lon, Lat, Alt, Avst, Mast, Pre, Suh, Ate, Mate, Aws, Mws, Hig, Set, Pop, GDP</i>

(2)	$x_\alpha = \sin \alpha$	α is the slope value.	<i>Slo</i>
(3)	$x_\gamma = \frac{\gamma}{100}$	γ is the humidity value.	<i>Arh</i>

190 2.3 Research Method

191 This study provides a methodological framework for predicting the occurrence of forest fires in China,
192 as shown in Figure 2. First, all the forest fire correlation factors are selected by feature selection to obtain
193 the forest fire driving factors that have great influence on fires. These factors are then used as input data
194 of the forest fire prediction model, and machine learning models (ANNs, RBF neural networks, SVMs
195 and RFs) are applied to obtain corresponding results. Finally, the model accuracy is obtained through
196 evaluation indexes such as the AUC value. The noun abbreviations in the following articles are shown in
197 Table 4.



198

199

Fig. 2 Flowchart of the Chinese forest fire occurrence prediction

200

Table 4: Abbreviations list

Full name	Abbreviation
Artificial Neural Network	ANN
Radial Basis Function	RBF

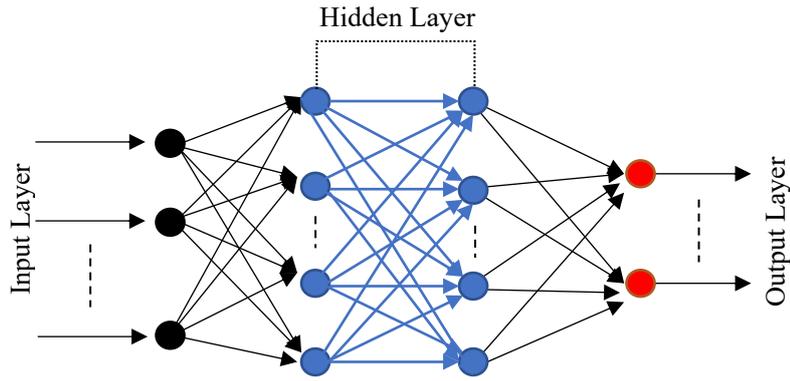
Radial Basis Function Neural Network	RBFNN
Support-Vector Machine	SVM
Random Forest	RF
Out-of-Bag	OOB
Area Under the Curve	AUC
True Positive	TP
True Negative	TN
False Positive	FP
False Negative	FN
Receiver Operating Characteristic	ROC
Support Vector Classification	SVC

201 2.3.1 Artificial Neural Networks

202 ANNs have become widely used in feedforward networks due to their clear structure, fast operation, easy
203 implementation, and abilities for self-learning and adaption to the environment[61][62]. ANNs consist
204 of three parts: an input layer, an output layer, and a hidden layer. The hidden layer may be a topological
205 structure of one or more layers, as shown in Figure 3. The input layer does not perform any calculations;
206 rather, it is used to receive data, that is, to transfer data to the adjacent hidden layer with different weights.
207 The hidden layer processes the data through a nonlinear activation function and then passes it to the
208 output layer. The final result is obtained from the output layer. The mathematical principle is as follows:

$$209 \begin{cases} h^{(1)} = \varphi^{(1)}(\sum_{i=1}^n x_i \cdot \omega_j^{(1)} + b^{(1)}) \\ y = \varphi^{(2)}(\sum_{j=1}^n h_i^{(1)} \cdot \omega_j^{(2)} + b^{(2)}) \end{cases} \quad (4)$$

210 In the formula, the input layer is $x \in R^m$, the hidden layer output is $h \in R^n$, the output layer is $y \in R^K$,
211 the input layer to the hidden layer weight connection matrix is $\omega^{(1)} \in R^{m \times n}$, the weight connection bias
212 from the input layer to the hidden layer is $b^{(1)} \in R^n$, and the weight connection matrix and the bias from
213 the hidden layer to the output layer are $\omega^{(2)} \in R^{n \times K}$ and $b^{(2)} \in R^{n \times K}$, respectively.



214

215

Fig. 3 Diagram of the structure of an ANN

216 2.3.2 Radial Basis Function Neural Network

217 The RBF neural network structure is a feedforward structure with an input layer, a single hidden layer,
 218 and an output layer. Its advantages are concise training and fast learning convergence speed, which can
 219 approximate any nonlinear function. This method has been widely used in time-series forecasting,
 220 nonlinear control systems, and the graphics-processing field. The basic idea of an RBF neural network
 221 is as follows. The RBF is used as the “base” of the hidden unit to form the hidden layer space. The hidden
 222 layer transforms the input vector and transforms the low-dimensional pattern input data into the high-
 223 dimensional space. The result is that the data are linearly separable in the high-dimensional space. The
 224 output of the RBF neural network is as follows:

225
$$y_i = \sum_{i=1}^h \omega_{ij} \exp\left(-\frac{1}{2\sigma^2} \|x_p - c_i\|^2\right) \quad j = 1, 2, \dots, n \quad (5)$$

226 where $x_p = (x_1^p, x_2^p, \dots, x_m^p)^T$ is the p^{th} input sample ($p = 1, 2, 3, \dots, P$), P is the total number of
 227 samples, c_i is the centre of the hidden layer node of the network, ω_{ij} is the connection weight from
 228 the hidden layer to the output layer, $i = 1, 2, 3, \dots, h$ is the number of hidden layer nodes, and y_i is the
 229 actual output of the j^{th} output node of the network corresponding to the input sample[63].

230 2.3.3 Support-Vector Machines

231 SVMs are mainly used for pattern classification and nonlinear regression. They are general learning
232 algorithms based on the principle of structural risk minimization. The core idea of an SVM is to establish
233 a classification hyperplane as a decision surface to maximize the isolation edge between the positive and
234 negative examples, thereby providing a high generalization performance[64]. SVMs can improve the
235 ability to transform data from high-dimensional spaces by flexibly using kernel functions when dealing
236 with various nonlinear problems. Taking a two-class SVM as an example, given a training set $T =$
237 $\{(x_1, y_1), \dots, (x_l, y_l)\} \in (X \times Y)^l$, where $x_i \in X = R^n, y_i \in \{1, -1\} (i = 1, 2, \dots, l)$, x_i is the feature
238 vector. The penalty parameter C and the kernel function $K(x, x')$ are first selected, and the optimization
239 problem is then constructed and solved as follows[63]:

$$240 \quad \min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l y_i y_j \alpha_i \alpha_j K(x_i, x_j) - \sum_{j=1}^l \alpha_j \quad (6)$$

$$241 \quad s. t. \quad \sum_{i=1}^l y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, i = 1, \dots, l \quad (7)$$

242 The optimal solution is then obtained: $\alpha^* = (\alpha_1^*, \dots, \alpha_l^*)^T$. A positive component of $\alpha^*: 0 \leq \alpha_j^* \leq C$
243 is then selected, and the threshold is calculated as follows:

$$244 \quad b^* = y_j - \sum_{i=1}^l y_i \alpha_i^* K(x_i - x_j) \quad (8)$$

245 Finally, the decision function is constructed:

$$246 \quad f(x) = \text{sgn}(\sum_{i=1}^l \alpha_i^* y_i K(x, x_i) + b^*) \quad (9)$$

247 2.3.4 Feature Selection and Random Forest

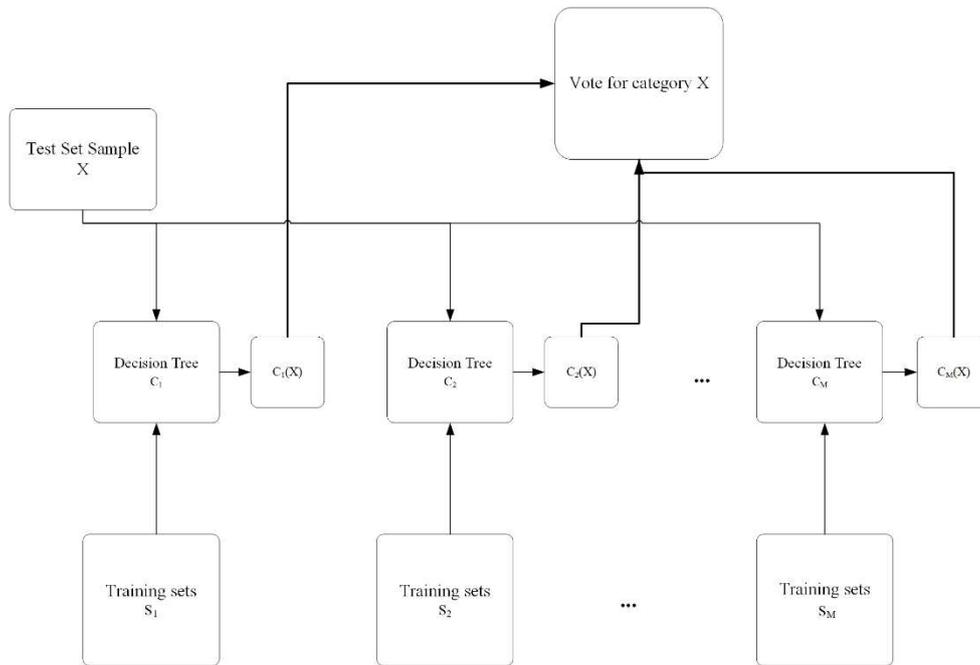
248 Feature selection refers to the selection of subsets from the original feature set to optimize a certain
249 evaluation criterion so that the model established with the optimal feature subset can achieve a prediction
250 accuracy similar to or better than that of the model established without feature selection[65][66]. RFs
251 have been demonstrated to have a high prediction accuracy and high tolerance to outliers and 'noise'[67].
252 This method can be used to evaluate the relationship between covariates and dependent variables and
253 calculate the relative importance of covariates[68][69]. RF has been applied in various fields in the past,
254 including medicine, genetics, ecology and remote sensing. In recent years, it has been widely used in
255 forest fire prediction and has demonstrated good predictive abilities[70][71][72][73]. The order of
256 importance of variables can be obtained by random forest algorithm. In previous studies, the variables
257 screened by this method have been proved to have high reliability[4][26][74][75][76]. Therefore, random
258 forest algorithm is selected as the method of feature selection in this study. The basic idea of feature
259 selection using RF is as follows: for the j th variable (X_j), the OOB error (err_{OOB}^j) of each tree t is
260 calculated, and then the value of the j th variable (X_j) is permuted while all others are left unchanged
261 among OOB data and the OOB error (err'_{OOB}^j) is again recalculated on this permuted dataset. RF
262 estimates the importance of a variable by evaluating how much the prediction error increases when the
263 OOB data for that variable are permuted. The importance score of X_j is as follows:

$$264 \quad VI(X^j) = \frac{1}{ntree} \sum_t (err'_{OOB}^j - err_{OOB}) \quad (10)$$

265 where R is the summation of all the trees, and $ntree$ is the number of trees in the RF[77][78]. For
266 classification, the OOB is the misclassification probability.

267 A RF is a highly flexible machine learning algorithm with broad application prospects. In essence, an RF

268 is a classifier consisting of multiple DTs formed by random methods. These trees are not related, hence
 269 its alternative name: “random decision tree.” When the test data enter the RF, each DT is classified, and
 270 the category with the most classification results among all the DTs is taken as the final result. The flow
 271 chart of the RF algorithm is shown in Figure 4.



272

273

Fig. 4 The flow chart of the RF algorithm

274 The basic principle of the RF algorithm is as follows. Let N be the number of attributes of the sample. n
 275 is an integer greater than 0 and less than N . First, the bootstrap method is used for resampling, randomly
 276 generating M training sets S_1, S_2, \dots, S_M . DTs A_1, A_2, \dots, A_M corresponding to each training set are then
 277 generated. Before selecting the attribute in each non-leaf node, n attributes are randomly selected from
 278 the N attributes as the split attribute set of the current node, and the node is split in the best split mode
 279 among the n attributes. Each tree grows intact without pruning. For the test set sample X , each DT is used
 280 to test and obtain the corresponding categories: $C_1(X), C_2(X), \dots, C_M(X)$. Finally, the voting method is
 281 adopted, and the category with the most output among the M DTs is regarded as the category to which

282 the test set sample X belongs[63].

283 **2.3.5 Model Performance Evaluation**

284 In this study, we used five performance indicators: accuracy, precision, recall, f1 value, and AUC to
285 evaluate the performance of the models. Descriptions of the five indicators are given below.

286 1. Accuracy: the proportion of the number of samples (TP and TN) that are correctly predicted to the total
287 number of samples. The formula is as follows:

$$288 \quad P = \frac{TP+TN}{TP+FP+TN+FN} \quad (11)$$

289 2. Precision: characterizes the classification effect of the classifier, which is the correct frequency value
290 predicted in the instance of the positive sample:

$$291 \quad T = \frac{TP}{TP+FP} \quad (12)$$

292 3. Recall: characterizes the recall effect of a certain class. It is the correct frequency of prediction in the
293 instance of the label as the positive sample:

$$294 \quad R = \frac{TP}{TP+FN} \quad (13)$$

295 4. f1 value: the value used to measure precision and recall. It is the harmonic mean of these two values:

$$296 \quad f1 = \frac{2TP}{2TP+FP+FN} \quad (14)$$

297 5. A receiver operating characteristic (ROC) curve is a method used to judge the prediction effect of the
298 model[64]. The prediction accuracy of the model is judged by the value of the AUC, which ranges from
299 0.5 to 1. The larger the value is, the closer the fit of the model is.

300 Note: TP, FN, FP, TN in the formulas are the labels of the confusion matrix form of the output result. The
301 form is shown in Table 5:

302 **Table 5: Confusion matrix form**

Prediction (column)/label (row)	Positive Sample	Negative Sample
Positive Sample	TP	FN
Negative Sample	FP	TN

303 **3. Results**

304 In this study, we used the MATLAB (MathWorks, USA, MATLAB 2019a) [61] and R Studio (JJ Allaire,
305 RStudio-1.2.5042/R 3.6.3) programming languages to implement the algorithms. We used MATLAB to
306 build the ANN, RBFNN, and SVM models and used R Studio to build the RF models.

307 To evaluate feature factors and model performance issues, the dataset was divided into two parts by
308 randomly selecting 70% of the pre-processed sample data as the training set and 30% as the test set [62]

309 **3.1 Feature Selection**

310 We use the RF algorithm to select the features of all variables after pre-treatment and select the subset of
311 features that had the greatest impact on the dependent variables for the next model construction process.

312 We divided the whole sample according to the above proportion (70% to the training set and 30% to the
313 test set) and repeated this process 5 times to obtain 5 training samples [4]. Then, the varSelRF package in
314 the R language was used to select and calculate the characteristic variables of the five training samples
315 to obtain the variable subsets of the five intermediate models, and the variables appearing more than
316 three times in the five variable quantum sets were selected as the variables after screening. The results

317 are shown in Table 6. All variables and variables filtered by feature selection were used as input data for
 318 RF modelling, and the out-of-pocket error rate (OOB) and confusion matrix were obtained as shown in
 319 Table 7 below.

320 **Table 6: Results of variable selection based on RF**

No.	Variable	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Frequency
1	Lat	+	+	+	+	+	5
2	Lon	+	+	+	+	+	5
3	Avst	+	+	+	+	+	5
4	Mast	+	+		+	+	4
5	Pre	+	+		+	+	4
6	Arh	+	+	+	+	+	5
7	Suh	+	+	+	+	+	5
8	Ate	+	+	+	+	+	5
9	Mate	+	+	+	+	+	5
10	Aws						0
11	Mws						0
12	Alt	+	+	+	+	+	5
13	Slo						0
14	Asp						0
15	Set						0
16	Hig						0
17	GDP	+	+		+	+	5
18	Pop	+	+	+	+	+	5

19	NDVI	+	+	+	+	+	5
20	Sfe						0

321

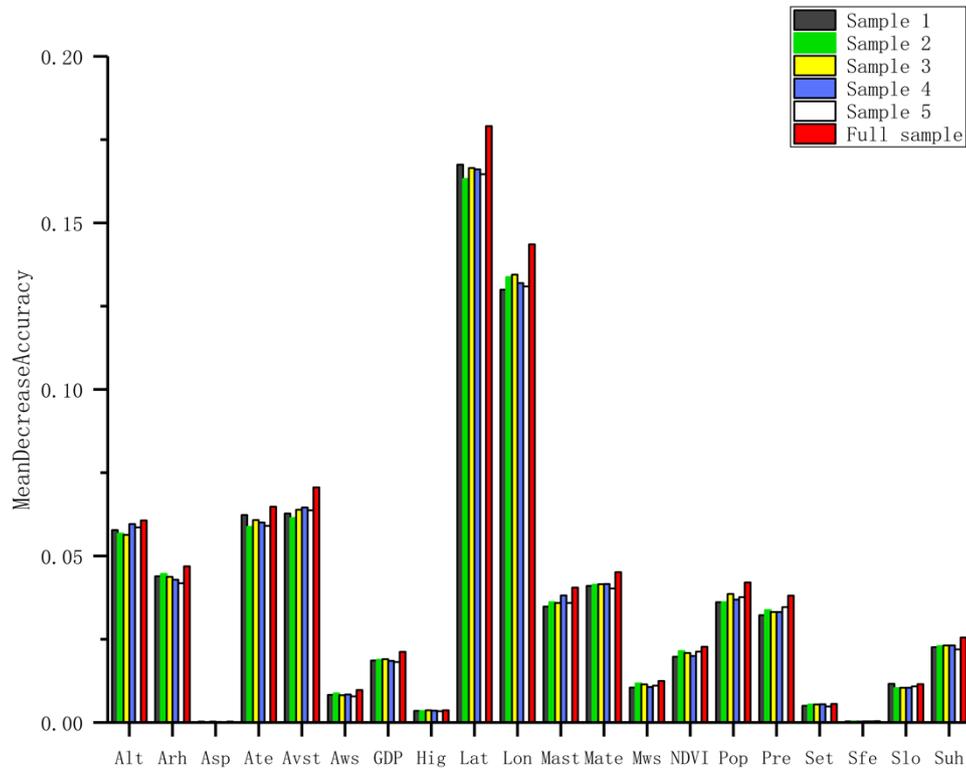
Table 7: The results of the OOB and confusion matrix of the two samples

Total variable sample	OOB estimate of error rate		10.89%
Confusion matrix:	0	1	Classification error rate
0	20224	2716	12.3%
1	2168	20737	9.5%
Sample of screened variables	OOB estimate of error rate		10.65%
Confusion matrix:	0	1	Classification error rate
0	20038	2810	11.8%
1	2171	20717	9.5%

322 It can be seen from Table 6 that the error rate outside the bag after using the whole variable modelling is
323 10.89%, while the error rate outside the bag after using the screened variable is 10.65%, which is lower
324 than the result of the whole variable. After feature filtering, the performance of the model is better, and
325 the complexity of the model is reduced, providing a simpler model. Finally, variables after feature
326 screening were taken as the main driving factors of forest fires and entered into the subsequent model
327 fitting process[26].

328 The results show that the main influencing variables are longitude, latitude, average surface temperature,
329 daily maximum surface temperature, accumulated precipitation, average relative humidity, sunshine
330 hours, average temperature, daily maximum temperature, altitude, population, GDP, and NDVI. These
331 variables performed subsequent model fitting. Then, the mean decrease in accuracy obtained by the RF
332 algorithm was used to evaluate the importance of the variable. The larger the value is, the greater the

333 importance of the variable is. Figure 5 shows the importance of each variable in the five random training
 334 samples and the 20 feature subsets in the full sample.



335

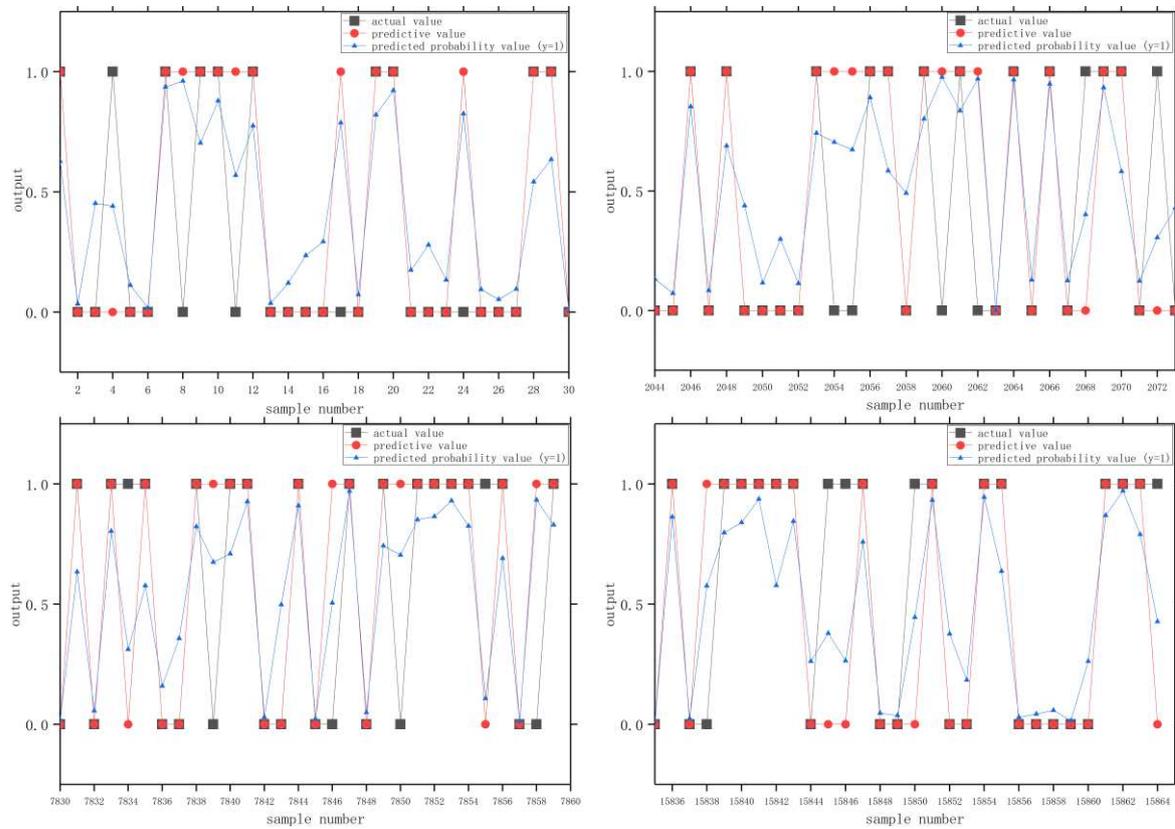
336 **Fig. 5 Feature subset importance**

337 **3.2 Model Fitting Results**

338 **3.2.1 Artificial Neural Network**

339 The input layer of the ANN consists of 13 neurons after feature selection: Lat, Lon, Avst, Mast, Pre, Arh,
 340 Suh, Ate, Mate, Alt, GDP, Pop, and NDVI. The output layer contains two cells (1 or 0). We use the
 341 gradient descent method to optimize the algorithm. We set the number of hidden layer cells between 1
 342 and 50, automatically select the optimal number of hidden layer cells as the final result, and finally obtain
 343 the number of hidden layer cells as 5. The comparison between the predictive value and the actual value
 344 in the test dataset is shown in Figure 6. Note: Due to the large sample size, only a part of the sample

345 comparison chart is displayed. This is also the case for the following comparison charts.



346

347

Fig. 6 Comparison charts of the predictive and actual values of the ANN

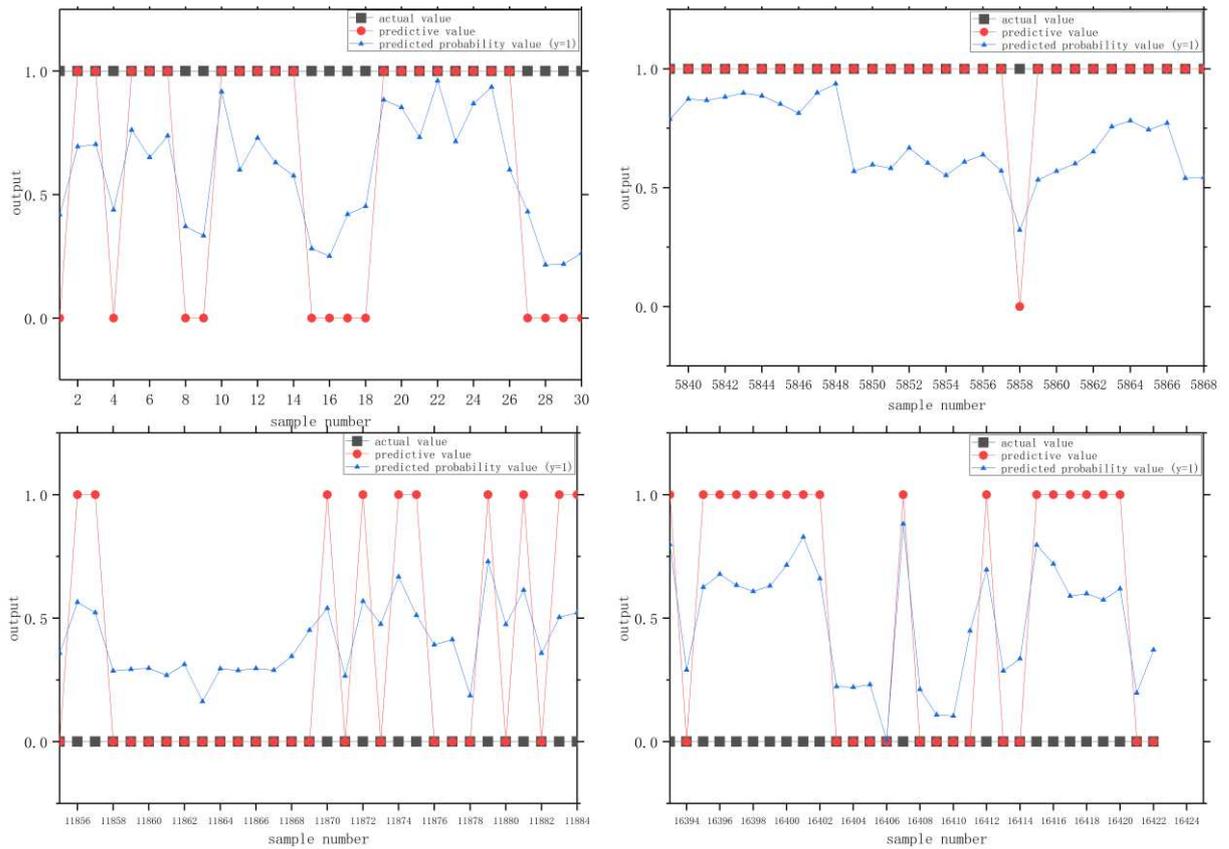
348 3.2.2 Radial Basis Function Neural Network

349 The input and output layer variables of the RBF neural network were the same as those of the ANN. The

350 number of hidden layers and the number of cells contained are the same as those in the MPNN model,

351 which are the optimal results automatically selected. After training, we obtained a hidden layer containing

352 10 units. The comparison charts of the predictive and actual values of the test set are shown in Figure 7.



353

354

Fig. 7 Comparison charts of the predictive and actual values of the RBFNN (part of the

355

sample)

356 3.2.3 Support-Vector Machine

357

We used the LIBSVM package of MATLAB to construct the SVM. The model was constructed using the

358

RBF kernel function for processing nonlinear data. We used the grid search method and 10-fold cross-

359

validation to select the parameters and determine the penalty parameter C and the kernel parameter g .

360

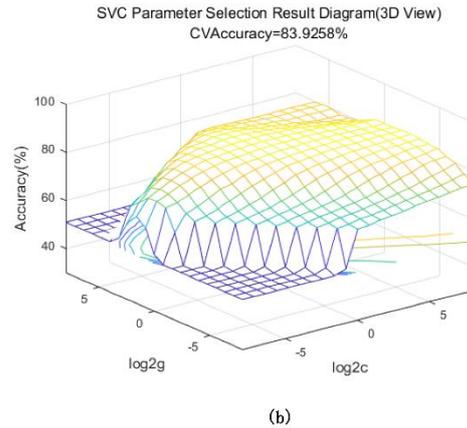
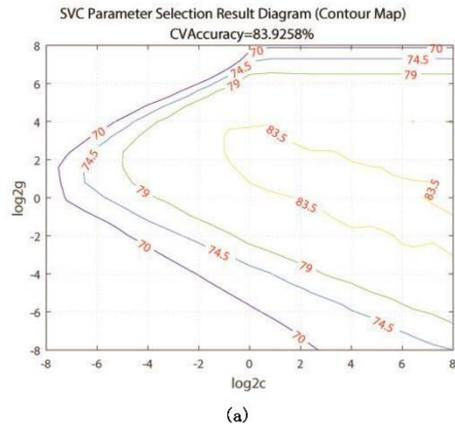
Figure 8 shows a contour map and a 3D view of the result of the SVC parameter selection. After

361

calculation, the accuracy rate of the grid search method reached 83.9%, and the accuracy rate of cross-

362

validation reached 82.6%.



363

364

Fig 8. SVC parameter selection result: (a) contour map (b) 3D view

365

It can be seen from the results that the optimal values of C and g are 1.74 and 3.03, respectively. After

366

setting the parameters to the optimal values, we performed SVM modelling and obtained the predicted

367

values. Figure 9 shows the comparison charts of the actual and predicted values. After optimization, the

368

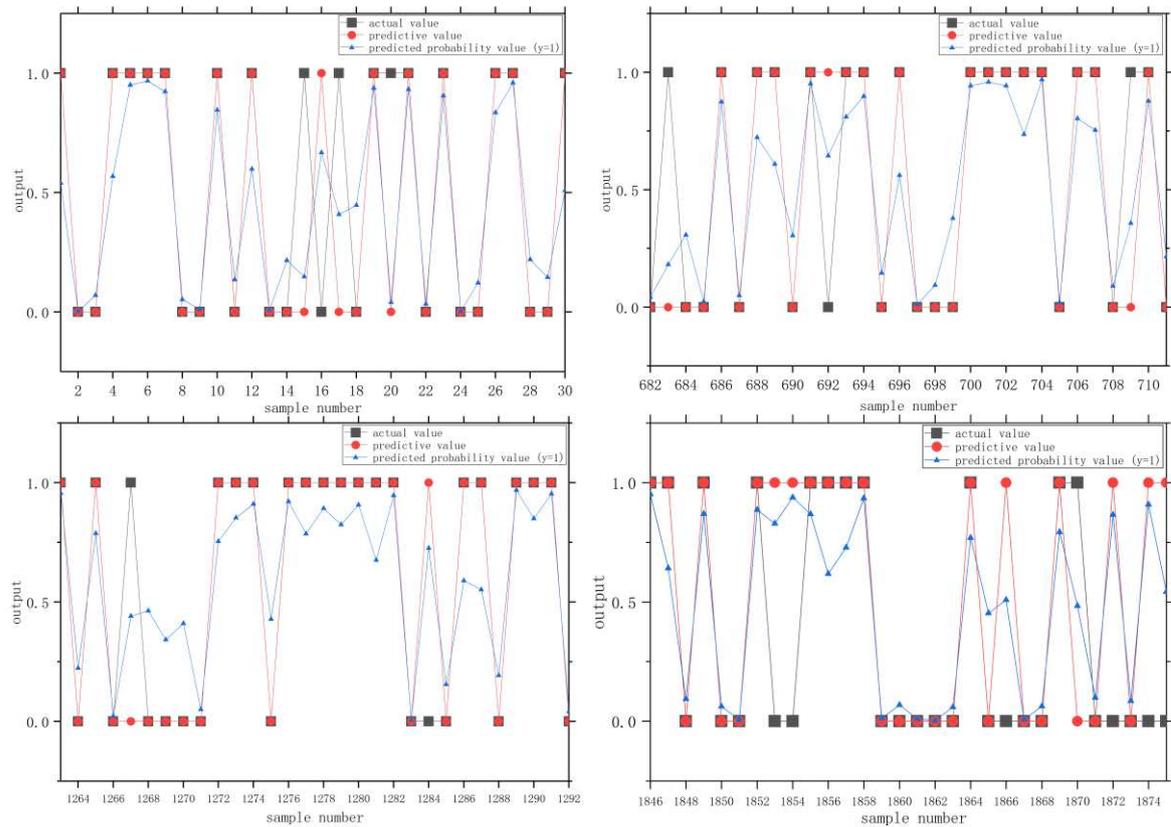
total number of support vectors was 19,460, and the number of support vectors at the boundary was

369

17,260. After model training, the accuracy rate of the training set was 86.02%, the accuracy rate of the

370

test set was 84.27%, and the performance of the model was high.



371

372

Fig. 9 Comparison charts of the predictive and actual values of the SVM (part of the sample)

373

374 3.2.4 Random Forest

375

We used the randomForest package in the R language to train the random training samples. We adjust

376

the ntree (number of DTs) and mtry (the node value of the trees) parameters. We then used cross-

377

validation to determine the optimal parameters of the model. Finally, we obtained the number of trees

378

and the accuracy of the test and training data through cross-validation. As shown in Figure 10, when the

379

number of DTs is 400 and the node value of the trees is 2, the accuracy tends to be stable. We used the

380

optimal number of DTs to create comparison charts of the actual and predicted values of the test set

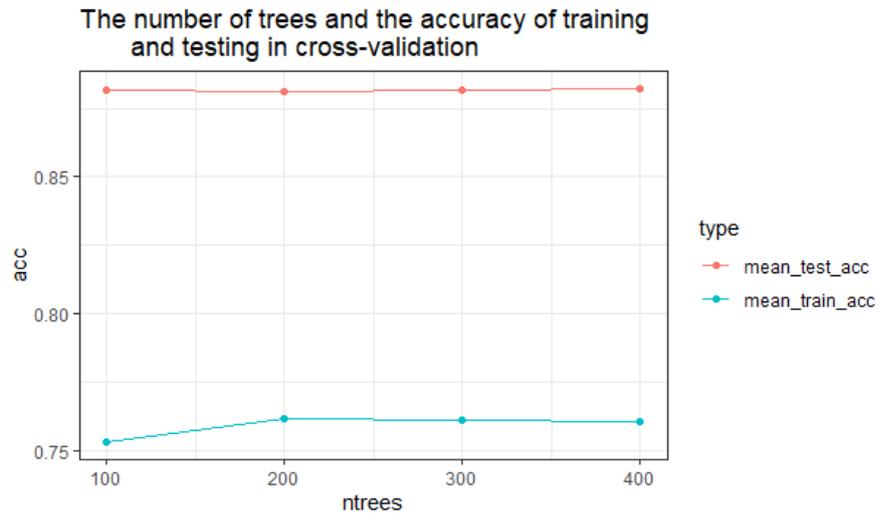
381

(Figure 11) and the average accuracy decline of 13 forest fire driving factors (Figure 12). Figure 10 shows

382

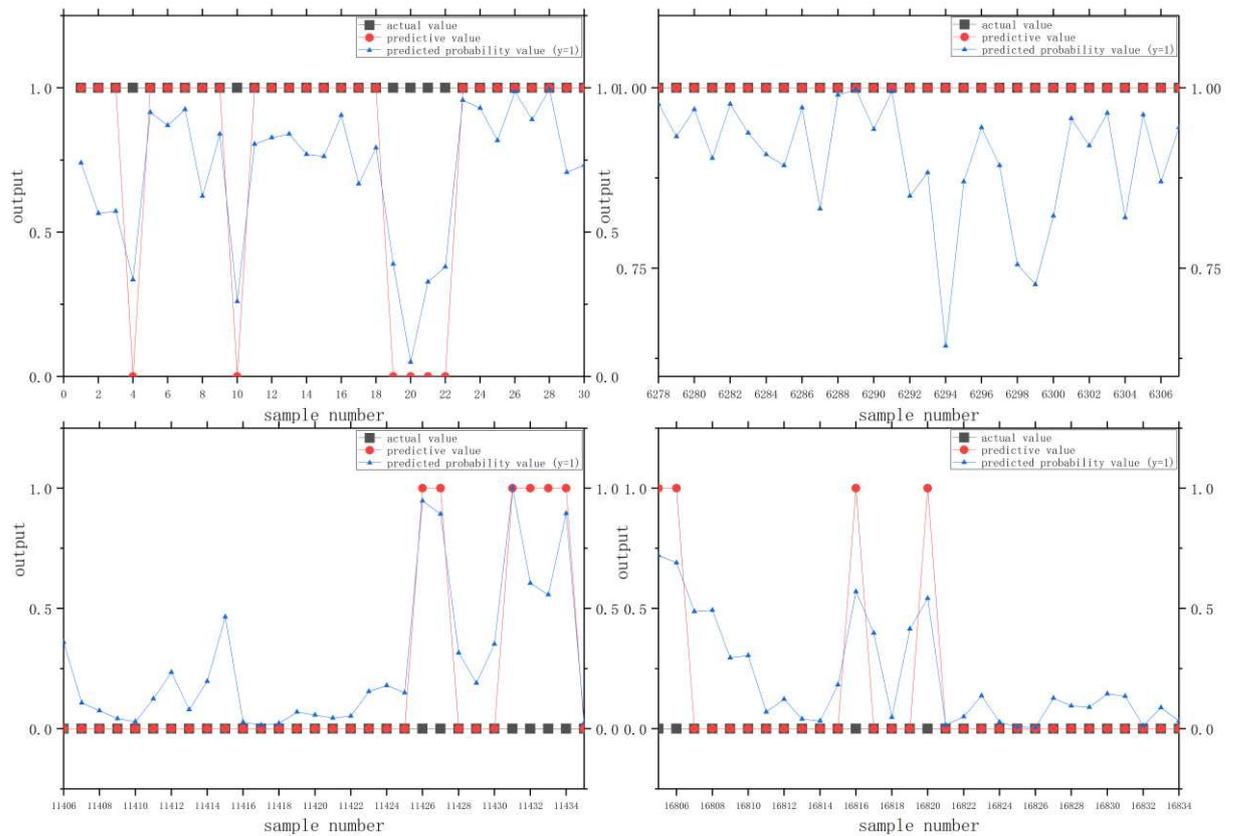
that among the main forest fire driving factors in China, the location variables that have the greatest

383 influence on the occurrence of forest fires are longitude and latitude. Rainfall is the variable with the
 384 smallest influence on the occurrence of forest fires.



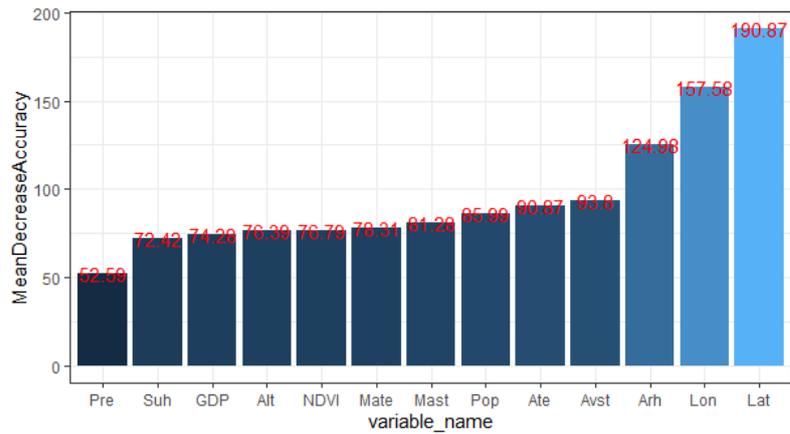
385

386 **Fig. 10** The number of trees and the accuracy of training and testing in cross-validation



387

388 **Fig. 11** Comparison charts of the predictive and actual values of the RF (part of the sample)



389

390

Fig. 12 Mean decrease accuracy of 13 variables

391 **3.3 Accuracy Evaluation**

392 We used the prediction results of the four models to construct a confusion matrix to obtain the accuracy,

393 precision, recall, f1 value, and AUC value, as shown in Table 8. Figure 11 shows the ROC curves of the

394 four models. As shown in Table 6, the accuracy and f1 values of each model were more than 75%, and

395 the AUC value was more than 0.80. Thus, the performance of all four models was high. Among the four

396 models, the RF model had the highest predictive ability, with an accuracy rate of 89.2%, a f1 value of

397 89%, and the highest AUC value, reaching 0.960. Compared with the other three models, the prediction

398 ability of the RBF neural network was the lowest, with an accuracy rate of 75.8% and an AUC value of

399 0.840. As shown in Figure 13, the RF model outperformed the other three models. We therefore

400 considered the RF model to be the most suitable among the four models for forest fire prediction in China.

401

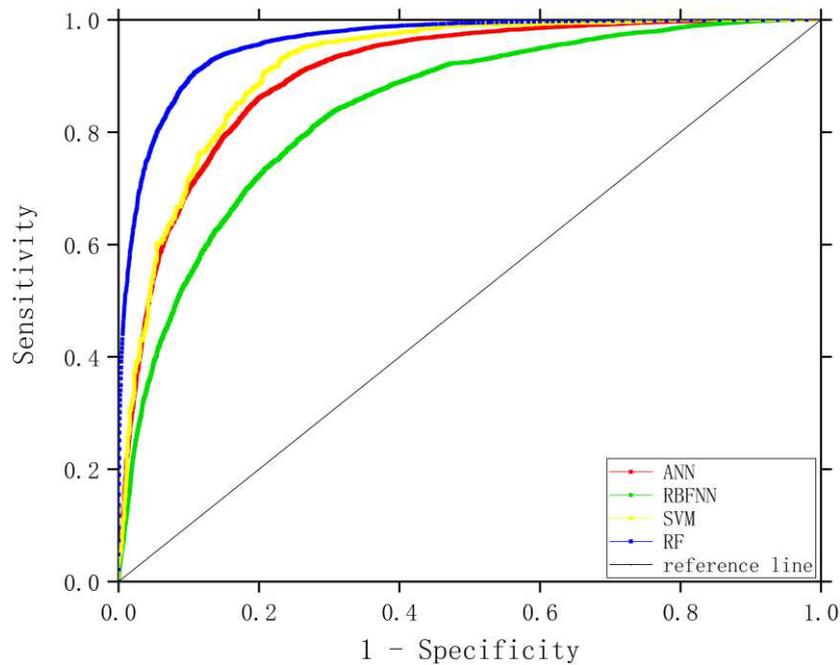
Table 8: Evaluation results of the four models

Model	Accuracy (%)	Precision (%)	Recall (%)	f1 value (%)	AUC
ANN	83.0	85.4	79.6	82.4	0.904

RBFNN	75.8	73.1	81.6	77.1	0.840
SVM	84.3	83.0	86.8	84.8	0.917
RF	89.2	90.2	87.9	89.0	0.960

402

Fig 11. Comparison charts of accuracy, precision, recall, and f1 values of the four models



403

Fig. 13 ROC curves of the four models

404

405 3.4 Fire Risk Classification

406 After evaluating the accuracy of the four models, we used the RF model (highest accuracy) to obtain the

407 probability of forest fire occurrence for the full sample. We used ArcGIS to draw a forest fire probability

408 map (Figure 14) and a seasonal forest fire probability map (Figure 15) for China. The numbers in the

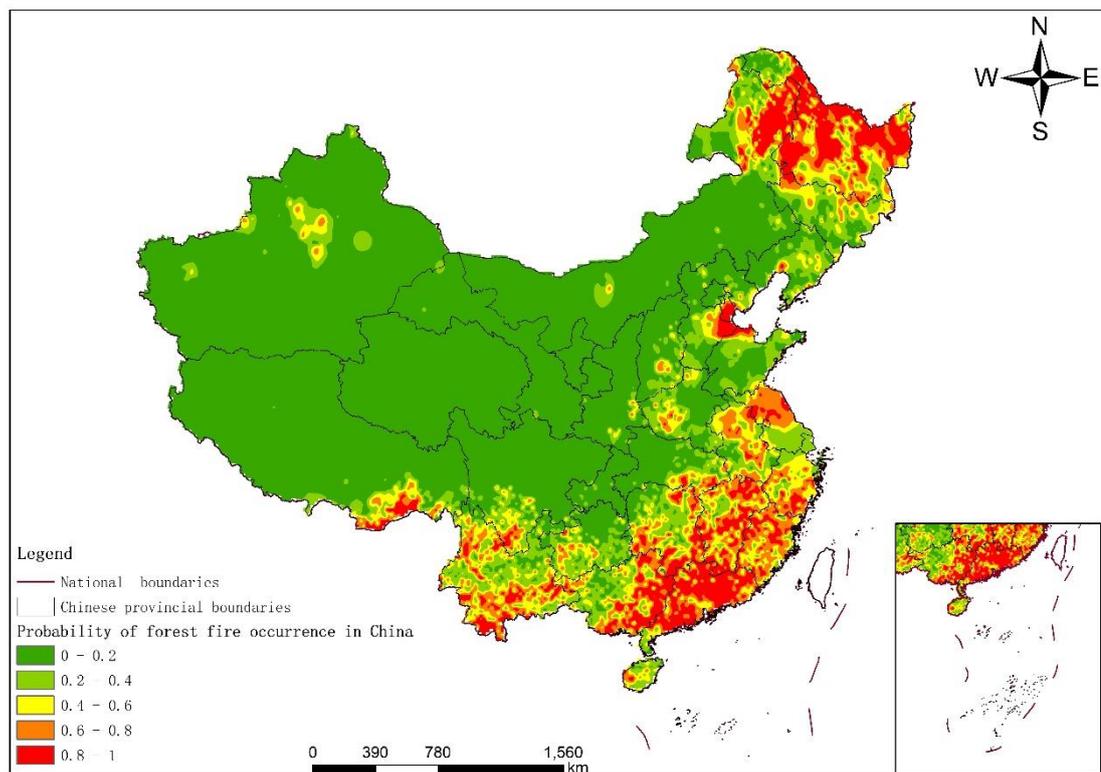
409 legends in Figure 14 and Figure 15 indicate the predicted value of the probability of forest fires in China.

410 For example, the probability of a forest fire is 1, which means that the probability of a forest fire is the

411 greatest; the number of red areas is 0.8-1, which indicates that the area is in a high-risk state, and forest

412 fires are very likely to occur. Figure 14 shows that the high incidence of forest fires in China is mainly

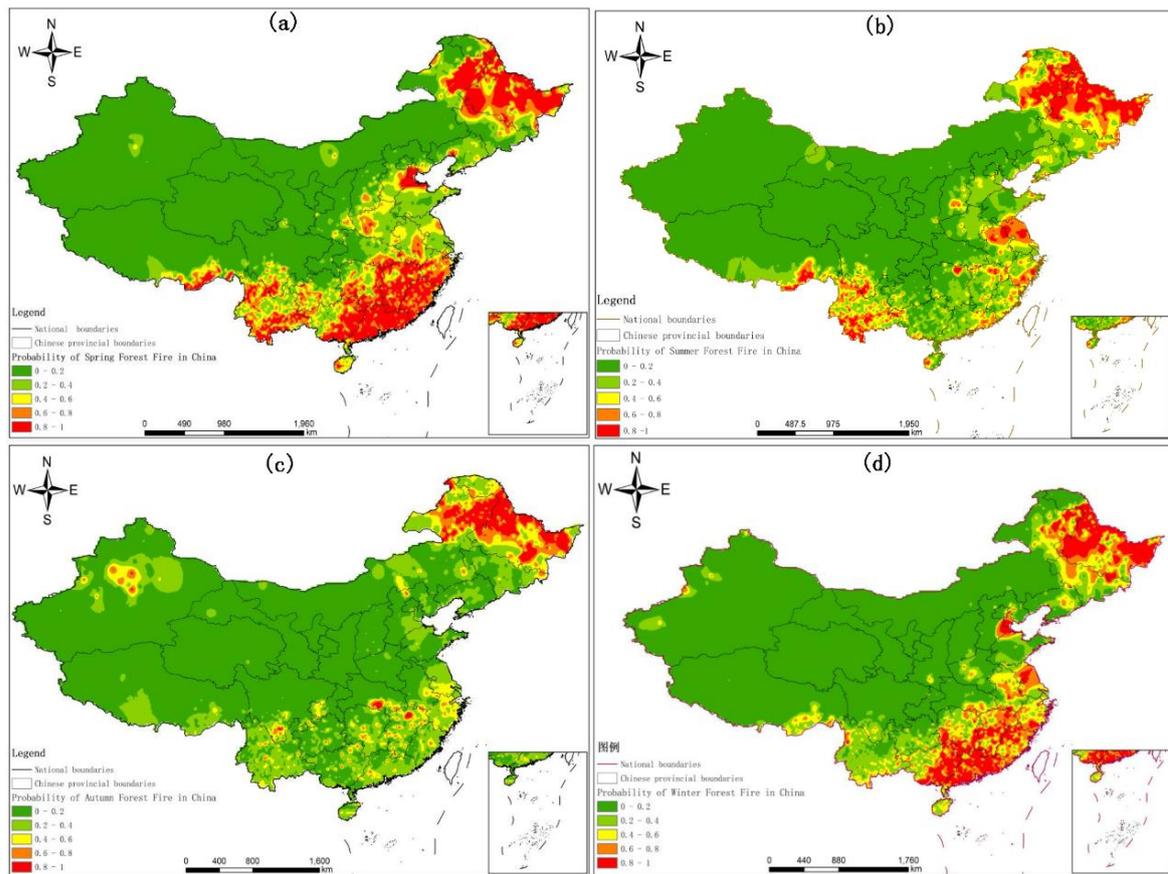
413 concentrated in the northeast (such as the Greater Xing'an Mountains region), the southeast (such as
414 Guangdong, Jiangxi, and Fujian), and the southwest (such as Yunnan and Sichuan). Overall, the
415 probability of forest fires in eastern China is higher than that in western regions, and the probability of
416 forest fires in the north and south is higher than that in Central China. Figure 15 shows that the seasonal
417 order of the probability of forest fires in China is, from highest to lowest, spring, winter, summer, and
418 autumn. Spring and winter are the seasons with a high incidence of forest fires, and the fires are mainly
419 concentrated in Northeast China (such as Heilongjiang Province) and south-eastern China (such as Fujian
420 Province and Guangdong Province).



421

422

Fig. 14 Forest fire probability map for China based on the RF model



423

424

Fig. 15 Seasonal forest fire probability map for China based on the RF model: (a) spring

425

(January, February, and March); (b) summer (April, May, and June); (c) autumn (July, August,

426

and September); and (d) winter (October, November, and December)

427

4. Discussion

428

4.1 Major Forest Fire Driving Factors in China and Their Impacts

429

In this study, 21 factors influencing the occurrence of forest fires were selected. These factors could be

430

divided into six categories: geographical location, meteorology, climate, topography, society and

431

vegetation. Researchers have studied the influencing factors of these forest fires[79][80][81][82].

432

Through feature selection, the driving factors of forest fires, such as longitude, latitude, average surface

433

temperature, daily maximum surface temperature, cumulative precipitation, average relative humidity,

434 sunshine duration, average temperature, daily maximum temperature, altitude, population, GDP and
435 NDVI, were obtained. In terms of the importance of feature subsets, longitude and latitude had the
436 greatest influence on the occurrence of forest fires. This result is because of the uneven distribution of
437 forest resources in China. Generally, in terms of forest resources, the south has more forested area than
438 the north, and the east has more forested area than the west. In addition, for example, many eucalyptus
439 trees are planted in south-eastern provinces, such as Guangdong and Fujian, which are prone to cause
440 fires. Therefore, the occurrence frequency of forest fires in China is obviously different with changes in
441 longitude and latitude. Therefore, longitude and latitude have become the most important factors in
442 studying the occurrence of forest fires in China. Currently, there are few studies regarding longitude and
443 latitude as the main driving factors of forest fires. However, some researchers have confirmed that the
444 number of forest fires increases with decreasing latitude [83][84]

445 Second, climatic factors have great impacts on forest fires, a result that is consistent with previous
446 research results[85][86]. Temperature is one of the three necessary conditions for combustion. When the
447 temperature reaches a certain level, forest fires are more likely to occur. The longer the daylight is, the
448 higher the temperature is and the greater the likelihood of forest fires is. Rainfall and average relative
449 humidity are other main factors affecting forest fires[56][87][88]. In addition, altitude and vegetation can
450 affect the occurrence of forest fires. Tian et al. (2013) believed that forest fires mainly occurred in low-
451 altitude areas [89]. Chuvieco et al. (2004) used the NDVI as a driving factor of forest fires to estimate
452 fuel humidity[90]. The higher the NDVI value is, the higher the vegetation coverage is. The more
453 flammable trees are, the more likely they are to cause problems associated with forest fires.

454 Forest fires are also driven by social and human factors (e.g., population and GDP). The larger the

455 population is, the greater the human activity level in the area is, and the greater the possibility of human-
456 caused forest fires is. Catry et al. (2007) and Sepulveda (2001) reached the same conclusion[91][92]. The
457 variables excluded in feature selection include the average wind speed, maximum wind speed, slope
458 direction, slope, nearest distance from the fire point to the highway, nearest distance from the fire point
459 to the residential area, and special festivals. This result indicated that in this experiment, these 7 variables
460 had little influence on the occurrence of forest fires. However, these factors may vary with time and space.
461 Other studies have found that these factors are the main driving factors of forest fires[93][94]. We believe
462 that this difference may be due to the different data selected and the different feature selection methods.
463 In future studies, multiple feature screening methods and analyses of different regions may be used to
464 obtain more comprehensive results.

465 **4.2 Optimal Choice of Forest Fire Prediction Model**

466 We entered the forest fire driving factors selected by feature selection into the four models (ANN,
467 RBFNN, SVM and RF) for training. We then evaluated them using five criteria: accuracy, precision,
468 recall rate, f1 value, and AUC value. We selected the RF model as the optimal choice for forest fire
469 prediction. The accuracies of all four models were above 75%, which meant that they were all reliable.
470 The RF model, however, exhibited the greatest prediction ability. The RBFNN model had the lowest
471 prediction performance.

472 Samaher et al. (2018)used a cascade correlation network, multilayer perceptron neural network,
473 polynomial neural network, RBF, and SVM for forest fire prediction [27]. They found that the prediction
474 performance of the SVM was the highest, and the performance of RBFNN was the lowest, which was
475 consistent with our conclusion. Sakr et al. (2011)used an SVM and an ANN to predict fire risk in Lebanon

476 [43]. Their results showed that the performance of the SVM model was higher than that of the ANN
477 model. This finding was similar to ours. Paulo et al. (2007) used RF and SVM models to predict forest
478 fires [8]. They concluded that the performance of the SVM model was higher than that of the RF model.
479 Their findings differed from our results. A possible reason for this difference is that Paulo et al. chose
480 four types of weather factors and made predictions for small areas, whereas we chose 13 types of factors
481 and made predictions for a large area. The choice of variables and the difference in the sample size affect
482 the model training.

483 Bisquert et al. (2012) used an ANN to establish a forest fire hazard model with the highest accuracy rate
484 of 76%, which was lower than the accuracy of our model (83%) [95]. Hong (2018) used an SVM algorithm
485 to analyse Dayu County in southwestern Jiangxi Province, China [96]. The results showed that the AUC
486 value of the SVM was 0.75, which was lower than the value in our model (0.92). Pourtaghid et al.
487 (2016) used an RF to conduct forest fire sensitivity analysis with a prediction accuracy of 72.8% [51].
488 Our model reached a prediction accuracy of 89.2%.

489 The four models we selected all exhibited high predictive capabilities. The main reason for this result
490 may be that appropriate multidimensional variables have been screened out and the data sample size is
491 large, which makes the training of each model more accurate and reliable.

492 There are also differences in the characteristics of these four models. The ANN and RBFNN models can
493 be trained very quickly, and they can handle samples with a large amount of data, but their accuracy in
494 this experiment is relatively low. In subsequent studies, particle swarms or genetic algorithms could be
495 used to improve the accuracy of these models. The SVM model has a high predictive ability, but it also
496 has shortcomings. The higher the model complexity is, the lower the calculation speed is. It takes a longer

497 time in this model to obtain the optimal parameters when processing large amounts of sample data. We
498 will consider using other algorithms to optimize the SVM model in the future. The RF model exhibited
499 excellent expressive power in this experiment. It can quickly process large data samples while ensuring
500 a high prediction accuracy.

501 **4.3 Recommendations for Forest Fire Prevention**

502 We produced a probability map for forest fires in China that showed that the highest incidences of forest
503 fires were in the northeast (Heilongjiang Province and the northern Inner Mongolia Autonomous Region),
504 the southeast (Fujian Province, Guangdong Province, and Jiangxi Province), and Yunnan Province. The
505 pattern of forest fire points presented a spatial clustering distribution. Ma et al. obtained similar results[4].
506 For these high-incidence areas, watch towers and monitoring equipment should be added for monitoring
507 and management. Moreover, the length of the forest fire barrier net should be increased to reduce the
508 spread of fires. In addition, the number of fire brigades and fire vehicles should be increased to enhance
509 the disaster-mitigation capabilities. Regarding seasonal forest fire risks, forest fire prevention and control
510 should be emphasized in spring and winter. Strengthening fire-prevention management during these
511 periods would mainly involve strengthening the management of human activities to reduce human-made
512 forest fires and improving publicity and education, such as the addition of fire-prevention signs.

513 This study has some shortcomings, and there is room for improvement. One of the three elements of fire
514 is combustible fuel. For the selection of forest fire driving factors, however, there is currently no way to
515 obtain data on fuel load and other related factors. Thus, this experiment lacked relevant data such as the
516 combustible load, particle size of combustible material, and combustible tree species. If possible, in
517 future research, such data could be added to the forest fire prediction model.

518 This study selected four kinds of machine learning algorithms for the forest fire prediction model. Other
519 applicable machine learning algorithms could be used in future experiments. In addition, the ability of
520 these machine learning algorithms to analyse spatial heterogeneity is relatively weak. Subsequent
521 research could use geographically weighted regression to build a high-precision forest fire prediction
522 model.

523 **5. Conclusion**

524 This study determined the main driving factors of forest fire occurrence in China through feature
525 selection. The main factors that affected the occurrence of forest fires to varying degrees were
526 meteorological, topographic, man-made, and vegetation factors. We built four forest fire prediction
527 models using the following machine learning algorithms: an ANN, an RBFNN, an SVM, and an RF. The
528 results of the evaluation showed that the accuracy of all the models was higher than 75%. Thus, these
529 models can be used to build forest fire prediction models. Among the four models, the RF model had the
530 highest comprehensive predictive ability, with an accuracy of 89.25%. It was therefore the optimal choice
531 for a forest fire prediction model in China.

532 We used the RF model to predict the probabilities of forest fires in China. Based on these probabilities,
533 we drew a map of the probability of forest fire occurrence in China and a map of the probability of forest
534 fires in China by season (spring, summer, autumn, and winter). Finally, based on these maps, we
535 identified the high-incidence areas and areas at risk of forest fires. We then put forward fire prevention
536 recommendations for the corresponding regions and seasons.

537 This research helps to understand the main forest fire driving factors in China. This study provides a
538 reference for the selection of high-precision forest fire prediction models. In addition, it provides

539 suggestions regarding the optimal times and locations of forest fire prevention in China. Finally, this
540 study provides guidance for China's forest fire prevention and control work.

541 **Data Availability**

542 The data used to support the findings of this study are available from the corresponding author upon
543 request.

544 **Conflicts of Interest**

545 The authors declare that they have no conflicts of interest.

546 **Funding Statement**

547 This research was jointly supported by the medium long-term project of "Precision Forestry Key
548 Technology and Equipment Research" (No. 2015ZCQ-LX-01) and the National Natural Science
549 Foundation of China (No. U1710123).

550 **Authors' contributions**

551 Yudong Li performed the experiment and wrote the manuscript;

552 Zhongke Feng contributed to the conception of the study;

553 Ziyu Zhao, Wenyuan Ma and Shilin Chen helped perform the analysis with constructive discussions;

554 Hanyue Zhang helped perform the data analysis.

555 **Acknowledgments**

556 We would like to acknowledge support from Beijing Key Laboratory for Precision Forestry, Beijing
557 Forestry University, as well as all the people who have contributed to this paper.

558 **References**

- 559 [1] Bhusal, Satish, Mandal, Ram (2020) Forest fire occurrence, distribution and future risks in
560 Arghakhanchi district, Nepal. Journal of Geography2(1):10-20.
561 <https://www.researchgate.net/publication/341701669>.
- 562 [2] Singh B. K., N. Kumar, P. Tiwari (2019) Extreme Learning Machine Approach for Prediction of
563 Forest Fires using Topographical and Metrological Data of Vietnam. 2019 Women Institute of
564 Technology Conference on Electrical and Computer Engineering (WITCON ECE)- Dehradun
565 Uttarakhand, India 104-112. DOI: 10.1109/WITCONECE48374.2019.9092926.
- 566 [3] D Mandallaz, R Ye (1997) Prediction of forest fires with Poisson models. Canadian Journal of Forest
567 Research 27(10): 1685-1694. DOI: 10.1139/x97-103.
- 568 [4] Ma W., Feng Z., Cheng Z., Chen S., Wang F. (2020) Identifying Forest Fire Driving Factors and
569 Related Impacts in China Using Random Forest Algorithm. Forests. 11(5), 507.
570 DOI:10.3390/f11050507.
- 571 [5] Dimopoulou Maria, Giannikos Ioannis (2004) Towards an integrated framework for forest fire
572 control. Eur. J. Oper. Res152(2):476-486. DOI: 10.1016/S0377-2217(03)00038-9.
- 573 [6] Flannigan MD, Krawchuk MA, Groot WJD, Wotton BM, et al. (2009) Implications of changing
574 climate for global wildland fire. International Journal of Wildland Fire18(5):483–507.DOI :
575 10.1071/WF08187.
- 576 [7] Xu ZQ, Su XY, Zhang Y (2012) Forest Fire Prediction Based on Support Vector Machine. Chinese
577 Agricultural Science Bulletin28(13):126-131. Available online:
578 http://en.cnki.com.cn/Article_en/CJFDTOTAL-ZNTB201213025.htm. (accessed on 24 February
579 2012).

- 580 [8] Paulo Cortez, Anibal Morais (2007) A Data Mining Approach to Predict Forest Fires using
581 Meteorological Data. *New Trends in Artificial Intelligence* 512-523. DOI: US7517766 B2.
- 582 [9] Naderpour M., Rizeei H. M., Khakzad N., Pradhan B. (2019) Forest fire induced Natech risk
583 assessment: a survey of geospatial technologies. *Reliability Engineering and System Safety*
584 191:106558. DOI: 10.1016/j.ress.2019.106558.
- 585 [10] Miguel Boubeta, María José, Lombardía, et al. (2015) Prediction of forest fires occurrences with
586 area-level poisson mixed models. *Journal of Environmental Management* 154:151-
587 158. DOI:10.1016/j.jenvman.2015.02.009.
- 588 [11] Kanga S, Kumar S, Singh SK(2017) Climate induced variation in forest fire using remote sensing
589 and GIS in Bilaspur district of Himachal Pradesh. *International Journal Of Engineering And*
590 *Computer Science* 6:21695–702. DOI: 10.18535/ijecs/v6i6.23.
- 591 [12] Erten E, Kurgun V, Musaoglu N(2004) Forest fire risk zone mapping from satellite imagery and
592 GIS: a case study. *XXth Congress of the International Society for Photogrammetry and Remote*
593 *Sensing, Istanbul*. pp: 222-230.
- 594 [13] Chuvieco E, Aguadoa I, Yebraa M(2010) Development of a framework for fire risk assessment
595 using remote sensing and geographic information system technologies. *Ecological Modelling*
596 221:46-58. DOI: 10.1016/j.ecolmodel.2008.11.017.
- 597 [14] Adab H, Devi Kanniah K, Solaimani K (2013) Modeling forest fire risk in the northeast of Iran
598 using remote sensing and GIS techniques. *Nat Hazards* 65(3):1723-1743. DOI: 10.1007/s11069-
599 012-0450-8.
- 600 [15] Sachdeva S, Bhatia T, Verma A. GIS-based evolutionary optimized gradient boosted decision trees
601 for forest fire susceptibility mapping. *Nat Hazards* 2018(92):1399–418. DOI:10.1007/s11069-018-

602 3256-5.

603 [16] Dhall A., Dhasade A., Nalwade A., et al. (2020) A survey on systematic approaches in managing
604 forest fires. *Applied Geography* 121:102266. DOI: 10.1016/j.apgeog.2020.102266.

605 [17] Maffei C., Menenti M. (2019) Predicting forest fires burned area and rate of spread from pre-fire
606 multispectral satellite measurements. *ISPRS Journal of Photogrammetry and Remote Sensing*
607 158:263-278. DOI: 10.1016/j.isprsjprs.2019.10.013.

608 [18] Shang C, Wulder M A, Coops N C, et al. (2020) Spatially-Explicit Prediction of Wildfire Burn
609 Probability Using Remotely-Sensed and Ancillary Data. *Canadian Journal of Remote Sensing*
610 46(8):1-17. DOI: 10.1080/07038992.2020.1788385.

611 [19] SU ZW, LIU AQ, et al. (2016) Driving factors and spatial distribution pattenen of forest fire in
612 Fujian Province. *JOURNAL OF NATURAL DISASTERS*25(2):110-119. DOI:
613 10.13577/j.jnd.2016.0213.

614 [20] CV Garcia, PM Woodard, SJ Titus, et al. (1995) A Logit Model for Predicting the Daily Occurrence
615 of Human Caused Forest-Fires. *International Journal of Wildland Fire*5(2):101-111. DOI: DOI:
616 10.1071/WF9950101.

617 [21] Futao Guo, Selvara Selvalakshmi, Fangfang Lin, et al. (2016) Geospatial information on
618 geographical and human factors improved anthropogenic fire occurrence modeling in the Chinese
619 boreal forest. *Canadian Journal of Forest Research* 46(4): 582-594. DOI:10.1139/cjfr-2015-0373.

620 [22] K.Z. Liu, L.F. Shu, F.J. Zhao et al. (2017) Research on spatial distribution of forest fire based on
621 satellite hotspots data and forecasting model. *Journal of Forestry Engineering*2(04):128-133.

622 Available online: http://en.cnki.com.cn/Article_en/CJFDTOTAL-LKKF201704021.htm. (accessed
623 on April 2014).

- 624 [23] B.Q. LIAO, J. WEI et al. (2008) Logistic and ZIP Regression Model for Forest Fire Data. FIRE
625 SAFETY SCIENCE3:143-149. Available online: [http://en.cnki.com.cn/Article_en/CJFDTOTAL-
HZKX200803002.htm](http://en.cnki.com.cn/Article_en/CJFDTOTAL-
626 HZKX200803002.htm). (accessed on March 2008).
- 627 [24] D Mandallaz, R Ye (1997) Prediction of forest fires with Poisson models. Canadian Journal of Forest
628 Research27(10): 1685-1694. DOI:10.1139/x97-103.
- 629 [25] F.T. GUO, H.Q. HU et al. (2010) Relationship between forest lightning fire occurrence and weather
630 factors in Daxing'an Mountains based on negative binomial model and zero-inflated negative
631 binomial models. Chinese Journal of Plant Ecology21(01):159-164. Available online:
632 http://en.cnki.com.cn/Article_en/CJFDTOTAL-ZWSB201005014.htm. (accessed on May 2010).
- 633 [26] H.L. Liang, Y.R. Lin, G. Yang et al. (2016) Application of Random Forest Algorithm on the Forest
634 Fire Prediction in Tahe Area Based on Meteorological Factors. SCIENTIA SILVAE
635 SINICAE.52(01):89-98. Available online: [http://en.cnki.com.cn/Article_en/CJFDTotal-
LYKE201601011.htm](http://en.cnki.com.cn/Article_en/CJFDTotal-
636 LYKE201601011.htm). (accessed on January 2016).
- 637 [27] Samaher Al-Janabia et al. (2018) Assessing the suitability of soft computing approaches for forest
638 fires prediction. Applied Computing and Informatics14(2):214-224. DOI:
639 10.1016/j.aci.2017.09.006.
- 640 [28] Volkan Sevinca, Omer Kucukb, Merih Goltasc (2020) A Bayesian network model for prediction and
641 analysis of possible forest fire causes. Forest Ecology and Management457:117723. DOI:
642 10.1016/j.foreco.2019.117723.
- 643 [29] Artés, T., Cencerrado, A., Cortés, A., and Margalef, T (2017) Time aware genetic algorithm for
644 forest fire propagation prediction: exploiting multi-core platforms. Concurrency Computat.: Pract.
645 Exper. 29(9): 3837. DOI:10.1002/cpe.3837.

- 646 [30] D.T. Bui, Q.T. Bui, Q.P. Nguyen, B. Pradhan, H. Nampak, P.T. Trinh (2017) A hybrid artificial
647 intelligence approach using GIS-based neural-fuzzy inference system and particle swarm
648 optimization for forest fire susceptibility modeling at a tropical area, *Agric. Agricultural and Forest
649 Meteorology*. 233:32-44. DOI:10.1016/j.agrformet.2016.11.002.
- 650 [31] M. Denham, A. Cortés, T. Margalef, E. Luque (2008) Applying a dynamic data driven genetic
651 algorithm to improve forest fire spread prediction. *International Conference on Computational
652 Science*. Springer Berlin Heidelberg 36-45. DOI:10.1007/978-3-540-69389-5_6.
- 653 [32] H. Hong, S.A. Naghibi, M.M. Dashtpajardi, H.R. Pourghasemi, W. Chen (2017) A comparative
654 assessment between linear and quadratic discriminant analyses (LDA-QDA) with frequency ratio
655 and weights-of-evidence models for forest fire susceptibility mapping in China. *Arabian Journal of
656 Geosciences* 10 (7):167. DOI:10.1007/s12517-017-2905-4.
- 657 [33] A.M. Özbayog ̇lu, R. Bozer (2012) Estimation of the burned area in forest fires using computational
658 intelligence techniques, *Procedia Computer Science* 12:282–287. DOI:
659 10.1016/j.procs.2012.09.070.
- 660 [34] You Y, Lu C, Wang W, Tang C-K (2019) Relative CNN-RNN: learning relative atmospheric
661 visibility from images. *IEEE Trans Image Process* 28(1):45–55.
662 <https://doi.org/10.1109/TIP.2018.2857219>.
- 663 [35] Govil K, Welch ML, Ball JT, Pennypacker CR (2020) Preliminary results from a wildfire detection
664 system using deep learning on remote camera images. *Remote Sensing* 12(1):166.
665 DOI:10.3390/rs12010166.
- 666 [36] Camp A, Oliver C, Hessburg P, Everett R (1997) Predicting late-successional fire refugia pre-dating
667 European settlement in the Wenatchee Mountains. *Forest Ecology and Management* 95(1):63-77.

- 668 DOI:10.1016/S0378-1127(97)00006-6.
- 669 [37] Hr Pourghasemi, Beheshtirad M, Pradhan B (2016) A comparative assessment of prediction
670 capabilities of modified analytical hierarchy process (M-AHP) and Mamdani fuzzy logic models
671 using Netcad-GIS for forest fire susceptibility mapping. *Geomat Nat Hazards Risk* 7(2):861–885.
672 <https://doi.org/10.1080/19475705.2014.984247>.
- 673 [38] Li Z, Huang Y, Li X, Xu L (2020) Wildland Fire Burned Areas Prediction Using Long Short-Term
674 Memory Neural Network with Attention Mechanism. *Fire Technology*. DOI:10.1007/s10694-020-
675 01028-3.
- 676 [39] H. Soliman, K. Sudan and A. Mishra (2010) A smart forest-fire early detection sensory system:
677 Another approach of utilizing wireless sensor and neural networks. *SENSORS*, 2010 IEEE, Kona,
678 HI 1900-1904. DOI:10.1109/ICSENS.2010.5690033.
- 679 [40] Çetin Elmas, Yusuf Sönmez (2011) A data fusion framework with novel hybrid algorithm for multi-
680 agent Decision Support System for Forest Fire. *Expert Systems with Applications*.38(8):9225-9236.
681 <https://doi.org/10.1016/j.eswa.2011.01.125>.
- 682 [41] Onur Satir, Suha Berberoglu & Cenk Donmez (2016) Mapping regional forest fire probability using
683 artificial neural network model in a Mediterranean forest ecosystem, *Geomatics. Natural Hazards*
684 *and Risk*7:1645-1658. DOI:10.1080/19475705.2015.1084541.
- 685 [42] Maeda E. E., Formaggio A. R., Shimabukuro Y. E., et al. (2009) Predicting forest fire in the
686 Brazilian Amazon using MODIS imagery and artificial neural networks. *International Journal of*
687 *Applied Earth Observation and Geoinformation* 11(4):265–272. DOI:10.1016/j.jag.2009.03.003.
- 688 [43] Sakr G E, Elhajj I H, Mitri G (2011) Efficient forest fire occurrence prediction for developing
689 countries using two weather parameters. *Engineering Applications of Artificial Intelligence*24(5):

690 888-894. DOI:10.1016/j.engappai.2011.02.017.

691 [44] B.C. Ko, K.H. Cheong, J.Y. Nam (2009) Fire detection based on vision sensor and support vector
692 machines. *Fire Safety Journal* 44 (3): 322-329. DOI:10.1016/j.firesaf.2008.07.006.

693 [45] D.W. Xie, S.L. Shi (2014) Prediction for burned area of forest fires based on SVM model. *Applied*
694 *Mechanics and Materials*513(5):4084-4089. DOI:10.4028/www.scientific.net/AMM.513-517.4084.

695 [46] J. Zhao, Z. Zhang, S. Han, et al. (2011) SVM based forest fire detection using static and dynamic
696 features. *Computer Science and Information Systems*8(3): 821–841.
697 DOI:10.2298/CSIS101012030Z.

698 [47] Cutler DR, Edwards TC, Beard KH (2007) Random forests for classification in Ecology.
699 *Ecology*88(11):2783-2792. DOI:10.1890/07-0539.1.

700 [48] Anantha M. Prasad, Louis R. Iverson, Andy Liaw (2006) Newer classification and regression tree
701 techniques: Bagging and random forests for ecological prediction. *Ecosystems*9(2):181-
702 199.DOI:10.1007/s10021-005-0054-1.

703 [49] Rodrigucs M, De la Riva J (2014) An insight into machines learning algorithms to model
704 humarraused wildfire or currence. *Environmental Modelling & Software*57:192-
705 201.DOI:10.1016/j.envsoft.2014.03.003.

706 [50] LIANG HL, LIN YR, YANG G, et al. (2016) Application of random forest algorithm on the forest
707 fire prediction in Tahe area based on meteorological factors. *Scientla Silvae Sinicae*52(1):89-
708 98.DOI:10.11707/j.1001-7488.20160111.

709 [51] Pourtaghi Z.S., Pourghasemi H. R., Aretano R., Semeraro T. (2016) Investigation of general
710 indicators influencing on forest fire and its susceptibility modeling using different data mining
711 techniques. *Ecological Indicators*. 64:72–84. DOI:10.1016/j.ecolind.2015.12.030.

- 712 [52] Tian X., Zhao F., Shu L., et al. (2013) Distribution characteristics and the influence factors of forest
713 fires in China. *Forest Ecology and Management* 310:460–467. DOI:10.1016/j.foreco.2013.08.025.
- 714 [53] Chang Y., Zhu Z., Bu R., et al. (2015) Environmental controls on the characteristics of mean number
715 of forest fires and mean forest area burned (1987–2007) in China. *Forest Ecology and Management*
716 356:13–21. DOI:10.1016/j.foreco.2015.07.012.
- 717 [54] Zhong M., Fan W., Liu T., Li P. (2003) Statistical analysis on current status of China forest fire
718 safety. *Fire Safety Journal* 38(3):257–269. DOI: 10.1016/S0379-7112(02)00079-6.
- 719 [55] LU A.F. (2011) Study on the relationship among forest fire, temperature and precipitation and its
720 spatial-temporal variability in China. Institute of Geographic Sciences and Natural Resources
721 Research12:1396–1400. Available online: [http://en.cnki.com.cn/Article_en/CJFDTOTAL-
722 HNNT201109040.htm](http://en.cnki.com.cn/Article_en/CJFDTOTAL-HNNT201109040.htm). (accessed on September 2011).
- 723 [56] Ying. L., Han. J., Du. Y., Shen Z. (2018) Forest fire characteristics in China: Spatial patterns and
724 determinants with thresholds. *For. Ecol. Manag*424:345–354. DOI: 10.1016/j.foreco.2018.05.020.
- 725 [57] F.T. Guo, Z.W. Su, G.Y. Wang et al. (2017) Understanding fire drivers and relative impacts in
726 different Chinese forest ecosystems. *Science of the Total Environment*605: 411-425. DOI:
727 10.1016/j.scitotenv.2017.06.219.
- 728 [58] GUO F, WANG G, SU Z, et al. (2016) What drives forest fire in Fujian, China? Evidence from
729 logistic regression and random forests. *International Journal of Wildland Fire*25(5):505-519. DOI:
730 10.1071/WF15121.
- 731 [59] Chang Y, Zhu Z L, Bu R C, et al. (2013) Predicting fire occurrence patterns with logistic regression
732 in Heilongjiang Province, China. *Landscape Ecology*28(10):1989-2004.
733 <https://doi.org/10.1007/s10980-013-9935-4>.

- 734 [60] XU X L. (2018) Spatial distribution data set of quarterly vegetation index (NDVI) in China. Data
735 Registration and Publishing System of Resources and Environmental Science Data Center of
736 Chinese Academy of Sciences (<http://www.resdc.cn/DOI>),2018.DOI:10.12078/2018060601.
- 737 [61] A. Subasi, (2007) EEG signal classification using wavelet feature extraction and a mixture of expert
738 model. *Expert Systems with Applications*32 (4): 1084-1093.DOI: 10.1016/j.eswa.2006.02.005.
- 739 [62] D. Chen (2019) Prediction of Forest Fire Occurrence in Daxing'an Mountains Based on Logistic
740 Regression Model. *FOREST RESOURCES MANAGEMENT* (01):116-122. Available online:
741 http://en.cnki.com.cn/Article_en/CJFDTotol-LYZY201901018.htm. (accessed on January 2019).
- 742 [63] X.C. Wang et al. (2013) 43 Cases of MATLAB neural network analysis. Beijing, Bei hang
743 University Press.
- 744 [64] Yudong Li, Zhongke Feng , Shilin Chen, Ziyu Zhao, and Fengge Wang (2020) Application of the
745 Artificial Neural Network and Support Vector Machines in Forest Fire Prediction in the Guangxi
746 Autonomous Region, China. *Discrete Dynamics in Nature and Society* 2020:14.
747 <https://doi.org/10.1155/2020/5612650>.
- 748 [65] Ganteaume A, Camia A, Jappiot M, et al. (2013) A review of the main driving factors of forest fire
749 ignition over Europe. *Environmental Management* 51(3): 651-662.DOI:10.1007/s00267-012-9961-
750 z.
- 751 [66] Zhu Feng, Lu Liu (2020) Estimation of forest biomass in beijing (china) using multisource remote
752 sensing and forest inventory data. *Forests* 11(2): 163.DOI:10.3390/f11020163.
- 753 [67] Breiman (2001) Random forests. *Machine Lear* 45:15–32. DOI: 10.1023/A:1010933404324.
- 754 [68] Cutler DR, Edwards TJ, Beard KH, et al. (2007) Random forests for classification in ecology.
755 *Ecology* 88(11):2783–2792.DOI:10.1890/07-0539.1.

- 756 [69] Kubosova K, Brabec K, Jarkovsky J, Syrovatka V (2010) Selection of indicative taxa for river
757 habitats: a case study on benthic macroinvertebrates using indicator species analysis and the random
758 forest methods. *Hydrobiologia* 651: 101-114. DOI:10.1007/s10750-010-0280-1.
- 759 [70] Zhangwen Su, Haiqing Hu, Guangyu Wang, et al. (2018) Using GIS and Random Forests to identify
760 fire drivers in a forest city, Yichun, China. *Geomatics, Natural Hazards and Risk* 9(1): 1207-1229,
761 DOI:10.1080/19475705.2018.1505667.
- 762 [71] Oliveira S, Oehler F, San-Miguel-Ayanz J, Camia A, Pereira J. 2012. Modeling spatial patterns of
763 fire occurrence in Mediterranean Europe using multiple regression and random forest. *Forest Ecol*
764 *Manage.* 275:117-129. DOI: 10.1016/j.foreco.2012.03.003.
- 765 [72] Rodrigues M, Riva JDL (2014) An insight into machine-learning algorithms to model human-
766 caused wildfire occurrence. *Environ Model Softw* 57:192-201. DOI: 10.1016/j.envsoft.2014.03.003.
- 767 [73] Kane VR, Lutz JA, Alina Cansler C, et al. (2015) Water balance and topography predict fire and
768 forest structure patterns. *Forest Ecol Manag* 338:1-13. DOI: 10.1016/j.foreco.2014.10.038.].
- 769 [74] ZA Fang, XY B (2020). Improving land cover classification in an urbanized coastal area by random
770 forests: the role of variable selection. *Remote Sensing of Environment* 251:112105. DOI:
771 10.1016/j.rse.2020.112105.
- 772 [75] Epifanio, I. (2017) Intervention in prediction measure: a new approach to assessing variable
773 importance for random forests. *BMC Bioinform* 18:230. [https://doi.org/10.1186/s12859-017-1650-](https://doi.org/10.1186/s12859-017-1650-8)
774 8.
- 775 [76] Chan, C. W., Desiré Paelinckx (2008). Evaluation of random forest and adaboost tree-based
776 ensemble classification and spectral band selection for ecotope mapping using airborne
777 hyperspectral imagery. *Remote Sensing of Environment* 112(6):2999-3011. DOI :

778 10.1016/j.rse.2008.02.011.

779 [77] Zhangwen Su, Haiqing Hu, Guangyu Wang, et al. (2018) Using GIS and Random Forests to identify
780 fire drivers in a forest city, Yichun, China. *Geomatics, Natural Hazards and Risk* 9(1): 1207-1229,
781 DOI:10.1080/19475705.2018.1505667.

782 [78] Gromping U. (2009) Variable importance assessment in regression: linear regression versus random
783 forest. *The American Statistician* 63(4):308–319.DOI:10.1198/tast.2009.08199.

784 [79] Ganteaume A, Camia A, Jappiot M, et al. (2013) A review of the main driving factors of forest fire
785 ignition over Europe. *Environmental Management* 51(3): 651-662.DOI: 10.1007/s00267-012-
786 9961-z.

787 [80] Syphard AD, Radeloff VC, Keuler NS, et al. (2008) Predicting spatial patterns of fire on a southern
788 California landscape. *International Journal of Wildland Fire* 17:602–613.DOI: 10.1071/WF07087.

789 [81] Pew KL, Larsen CPS (2001) GIS analysis of spatial and temporal patterns of human-caused
790 wildfires in the temperate rainforest of Vancouver Island, Canada. *Forest Ecol Manag* 140:1-18.DOI:
791 10.1016/S0378-1127(00)00271-1.

792 [82] Dickson BG, Prather JW, Xu Y. et al. (2006) Mapping the probability of large fire occurrence in
793 northern Arizona, USA. *Landscape Ecol* 21: 747–761. <https://doi.org/10.1007/s10980-005-5475-x>.

794 [83] Ying L., Han J., Du Y., Shen Z. (2018) Forest fire characteristics in China: Spatial patterns and
795 determinants with thresholds. *Forest Ecology and Management* 424:345-
796 354.DOI:10.1016/j.foreco.2018.05.020.

797 [84] Prasad A.M., Iverson L.R., Liaw A. (2006) Newer classification and regression tree techniques:
798 Bagging and random forests for ecological prediction. *Ecosystems* 9:181-199.
799 DOI :10.1007/s10021-005-0054-1.

- 800 [85] Liu Z., Yang J., Chang Y., et al. (2012) Spatial patterns and drivers of fire occurrence and its future
801 trend under climate change in a boreal forest of northeast China. *Global Change Biology* 18:2041-
802 2056.DOI:10.1111/j.1365-2486.2012.02649.x.
- 803 [86] Syphard A. D., Radeloff V. C., Keuler N. S., et al. (2008) Predicting spatial patterns of fire on a
804 southern California landscape. *International Journal of Wildland Fire* 17:602-
805 613.DOI:10.1071/WF07087.
- 806 [87] Zumbrunnen T., Pezzatti G. B., Menéndez P., et al. (2011) Weather and human impacts on forest
807 fires: 100 years of fire history in two climatic regions of Switzerland. *Forest Ecology and*
808 *Management* 261:2188-2199. DOI:10.1016/j.foreco.2010.10.009.
- 809 [88] Cardille J A, Ventura S J, Turner M G, et al. (2001) ENVIRONMENTAL AND SOCIAL FACTORS
810 INFLUENCING WILDFIRES IN THE UPPER MIDWEST, UNITED STATES. *Ecological*
811 *Applications* 11(1): 111-127. DOI:10.1890/1051-0761(2001)011[0111: EASFIW]2.0.CO;2.
- 812 [89] Tian X., Zhao F., Shu L., Wang M. (2013) Distribution characteristics and the influence factors of
813 forest fires in China. *Forest Ecology and Management* 310:460-
814 467.DOI:10.1016/j.foreco.2013.08.025.
- 815 [90] Chuvieco E., Cocero D., Riaño D., et al. (2004) Combining ndvi and surface temperature for the
816 estimation of live fuel moisture content in forest fire danger rating. *Remote Sensing of Environment*
817 92: 322-331. DOI:10.1016/j.rse.2004.01.019.
- 818 [91] Catry F. X., Damasceno P., Silva J. S., et al. (2007). Spatial distribution patterns of wildfire ignitions
819 in Portugal. *Modelação espacial do risco de ignição em Portugal Continental*, 8. Project: Fire
820 ecology and post-fire restoration. Available online:
821 https://www.researchgate.net/publication/240613824_Spatial_Distribution_Patterns_of_Wildfire

822 [Ignitions in Portugal](#). (accessed on January 2007).

823 [92] Avila-Flores D., Pompa-Garcia M., Antonio-Nemiga X., et al. (2010). Driving factors for forest fire
824 occurrence in Durango State of Mexico: A geospatial perspective. *Chinese Geographical Science*
825 20(6), 491-497. DOI:10.1007/s11769-010-0437-x.

826 [93] Maingi K J, Henry M C (2007) Factors influencing wildfire occurrence and distribution in eastern
827 Kentucky, USA. *International Journal of Wildland Fire* 16(1): 23-33. DOI: 10.1071/ WF06007.

828 [94] Avilaflores D Y, Pompagarcia M, Antonionemiga X, et al. (2010) Driving factors for forest fire
829 occurrence in Durango State of Mexico: A geospatial perspective. *Chinese Geographical Science*.
830 20(6): 491-497. DOI:10.1007/s11769-010-0437-x.

831 [95] Bisquert M, Caselles E, Sanchez J M, et al. (2012) Application of artificial neural networks and
832 logistic regression to the prediction of forest fire danger in Galicia using MODIS data. *International*
833 *Journal of Wildland Fire* 21(8): 1025-1029. DOI:10.1071/WF11105.

834 [96] Hong, H., et al. (2018) Applying genetic algorithms to set the optimal combination of forest fire
835 related variables and model forest fire susceptibility based on data mining models. The case of Dayu
836 County, China. *Science of The Total Environment* 630:1044-1056. DOI:
837 10.1016/j.scitotenv.2018.02.278.

Figures

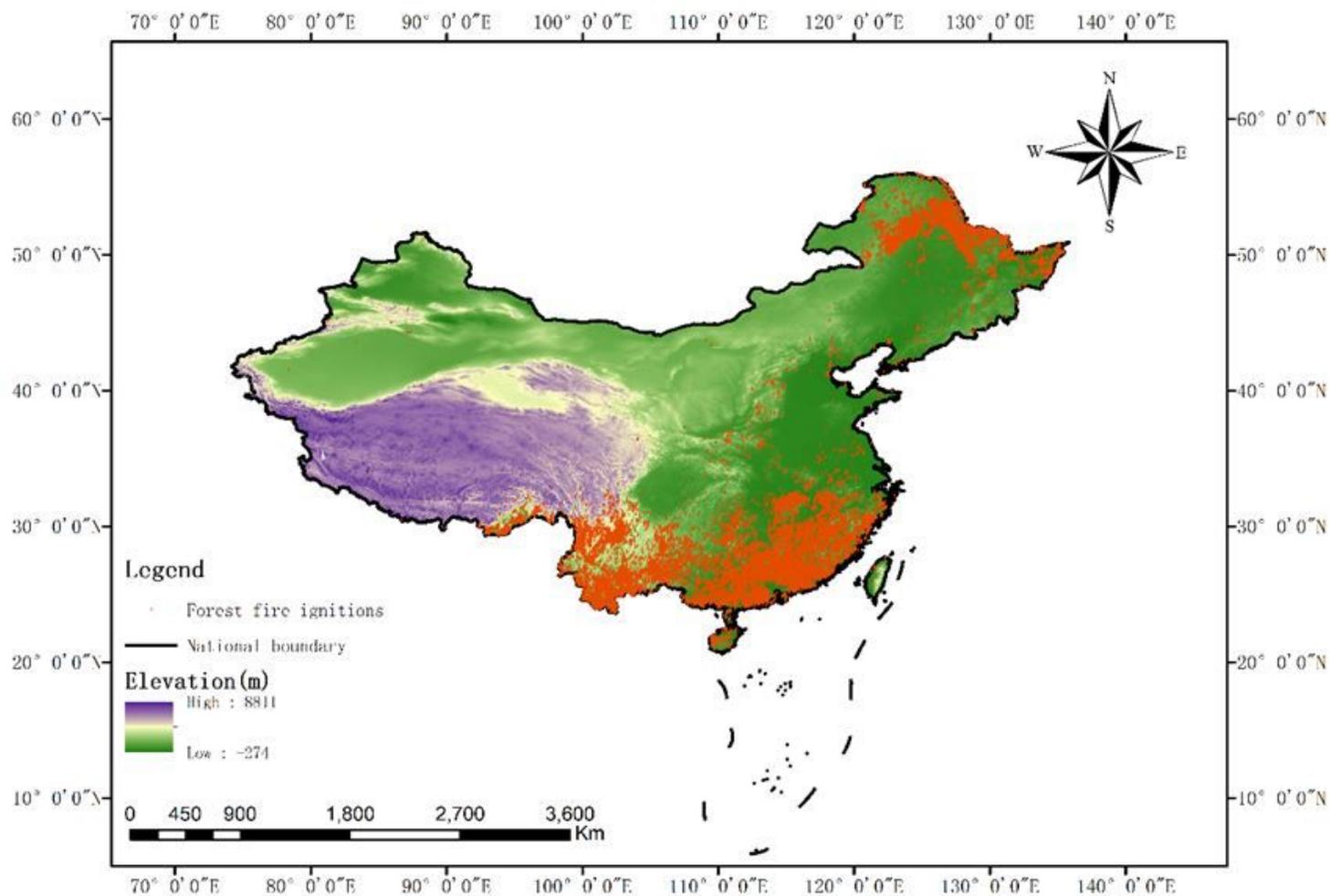


Figure 1

Map of the study area. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

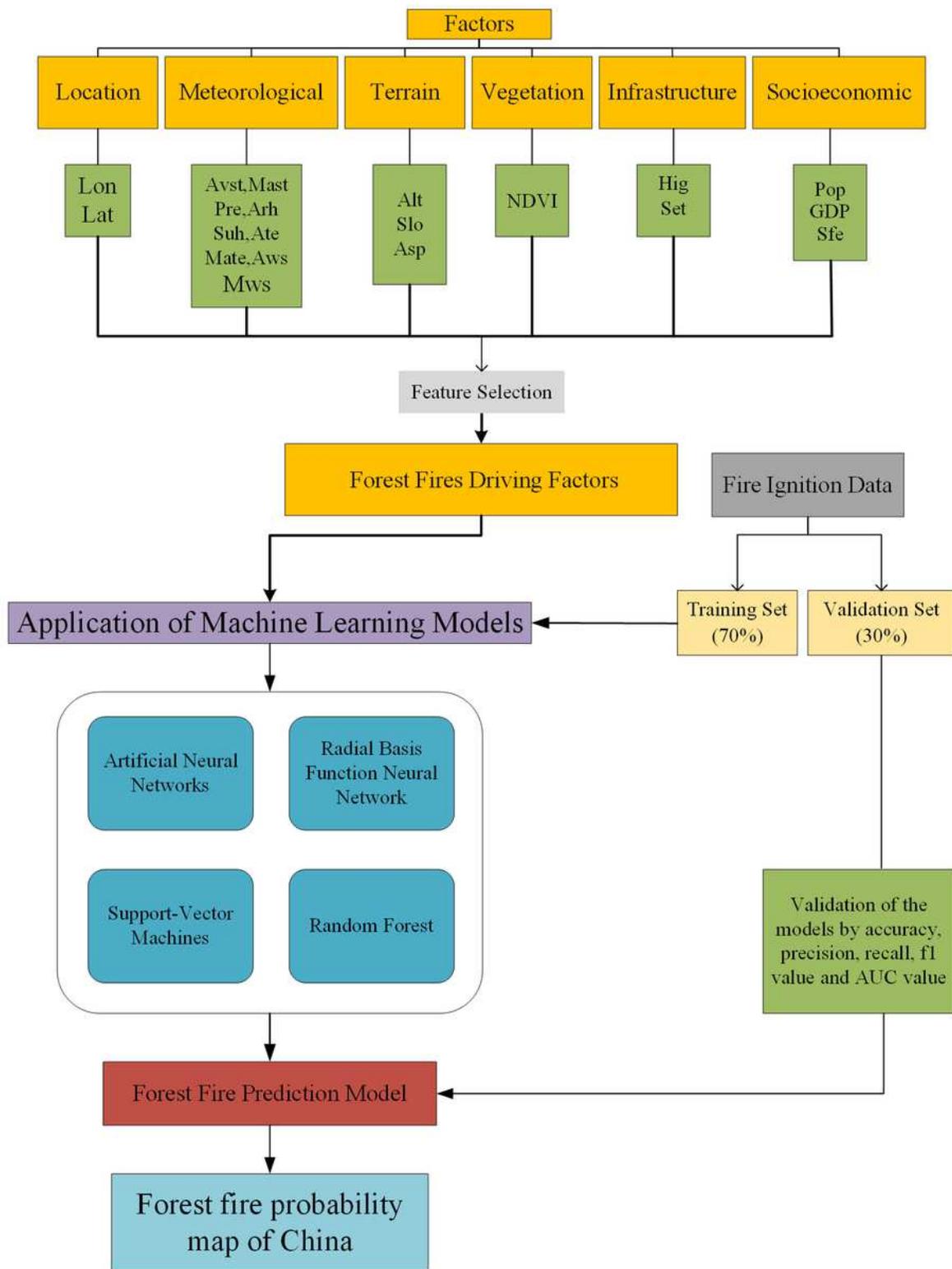


Figure 2

Flowchart of the Chinese forest fire occurrence prediction.

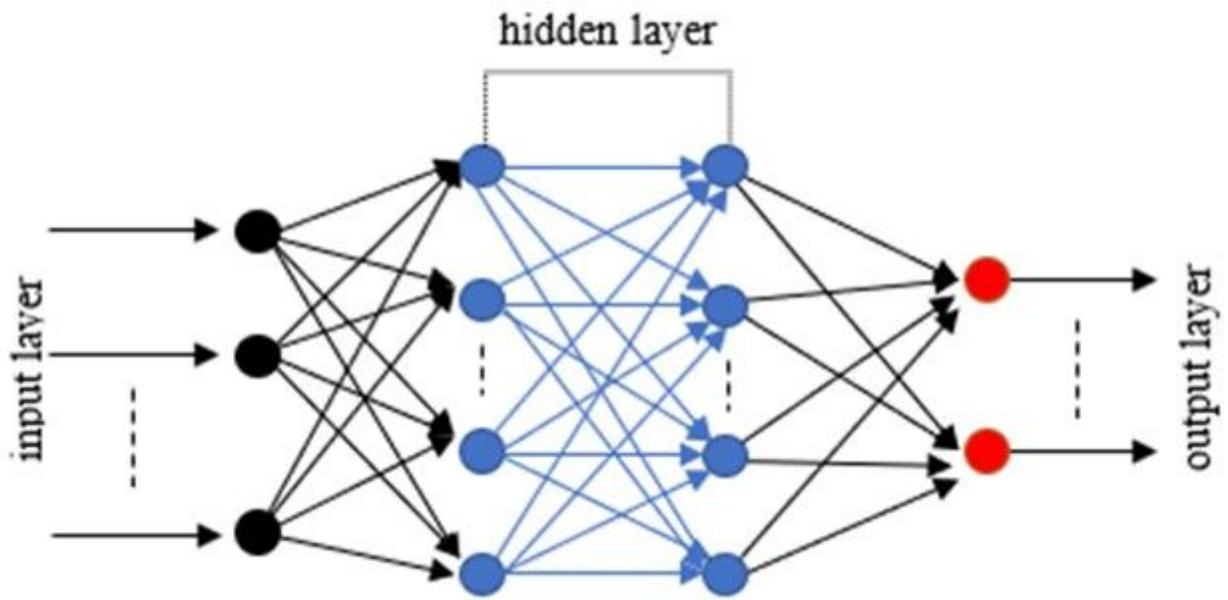


Figure 3

Diagram of the structure of an ANN.

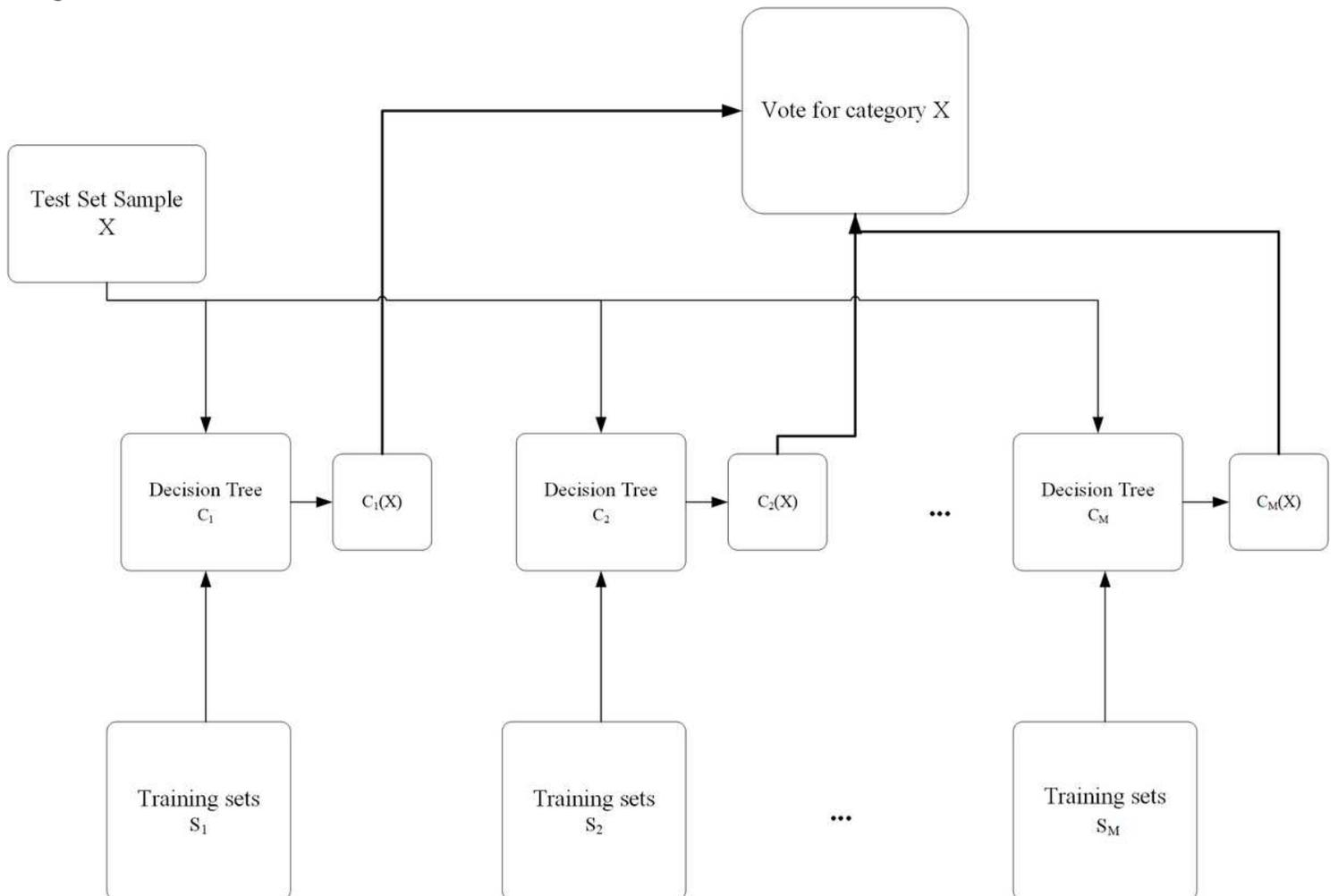


Figure 4

The flow chart of the RF algorithm

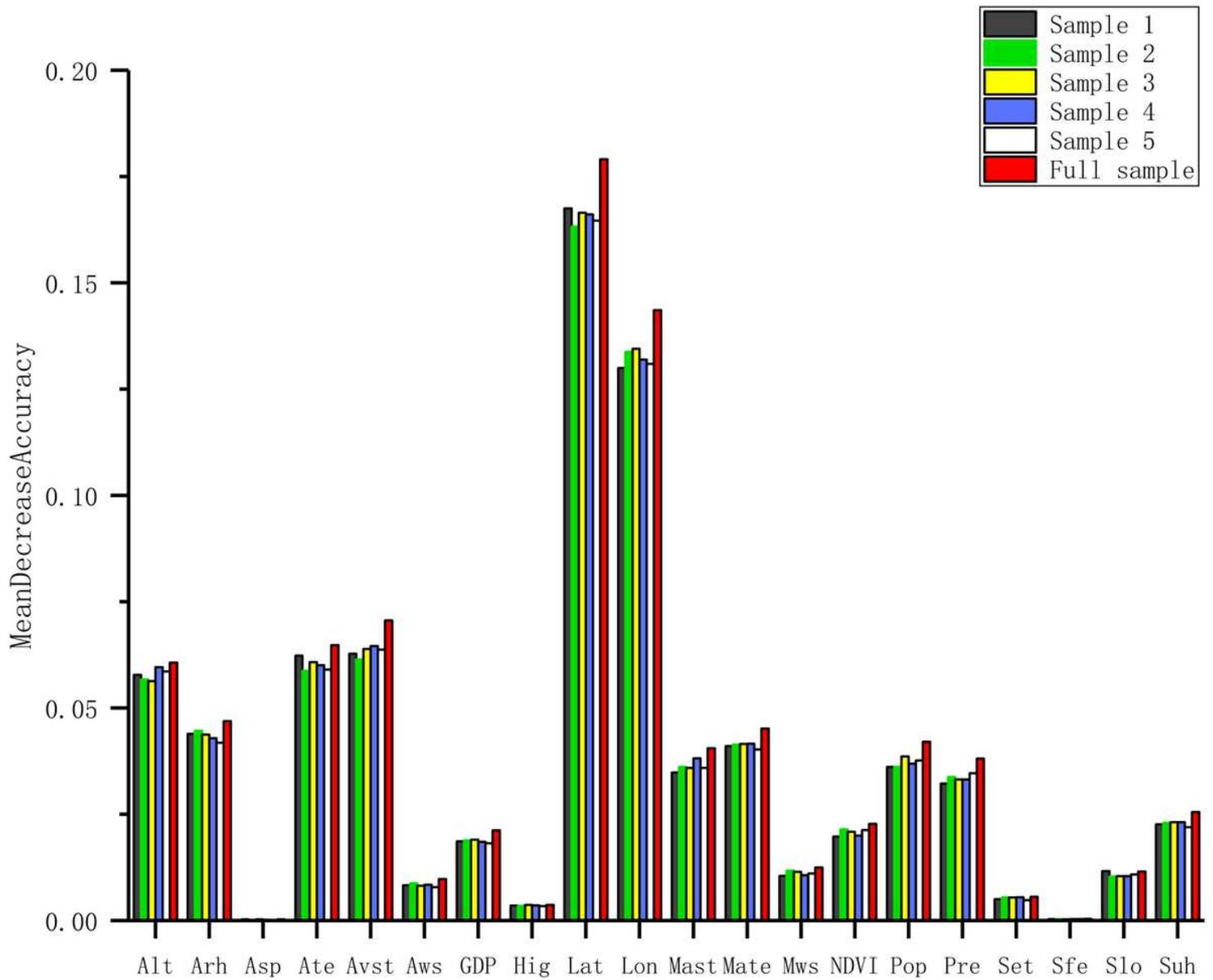


Figure 5

Feature subset importance.

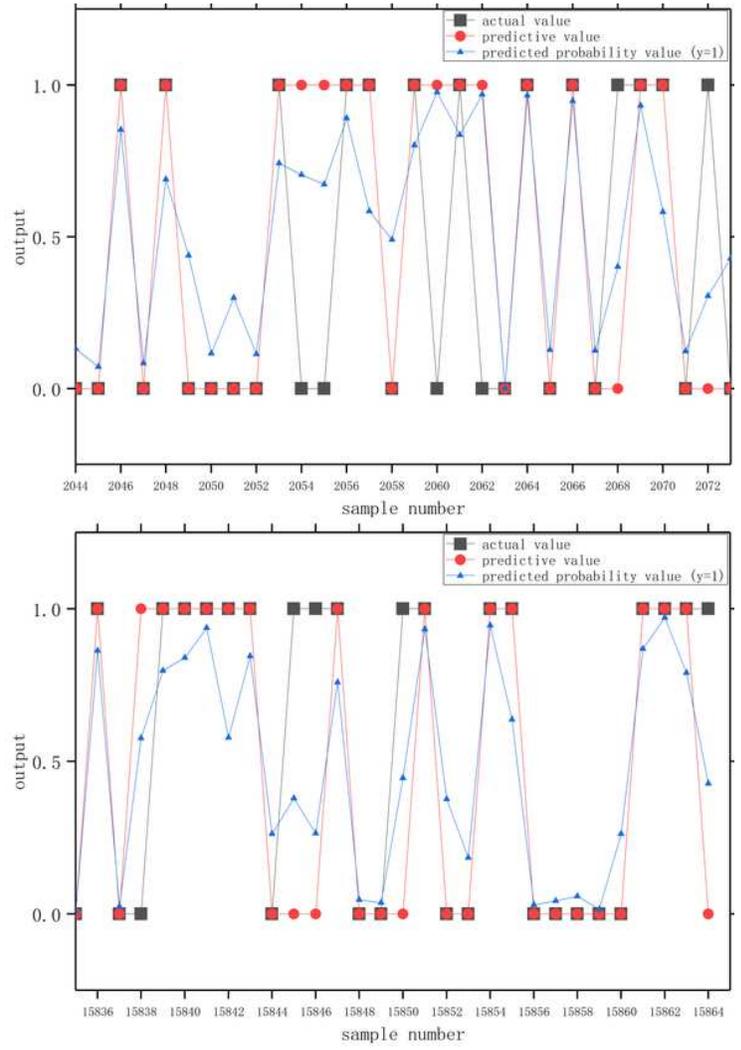
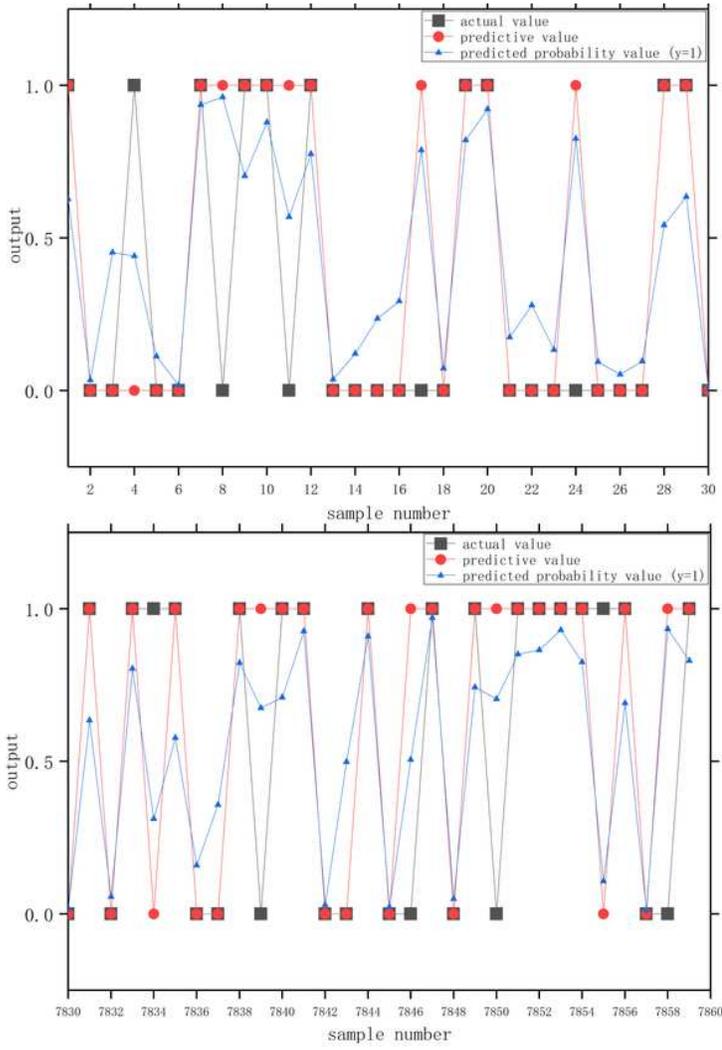


Figure 6

Comparison charts of the predictive and actual values of the ANN.

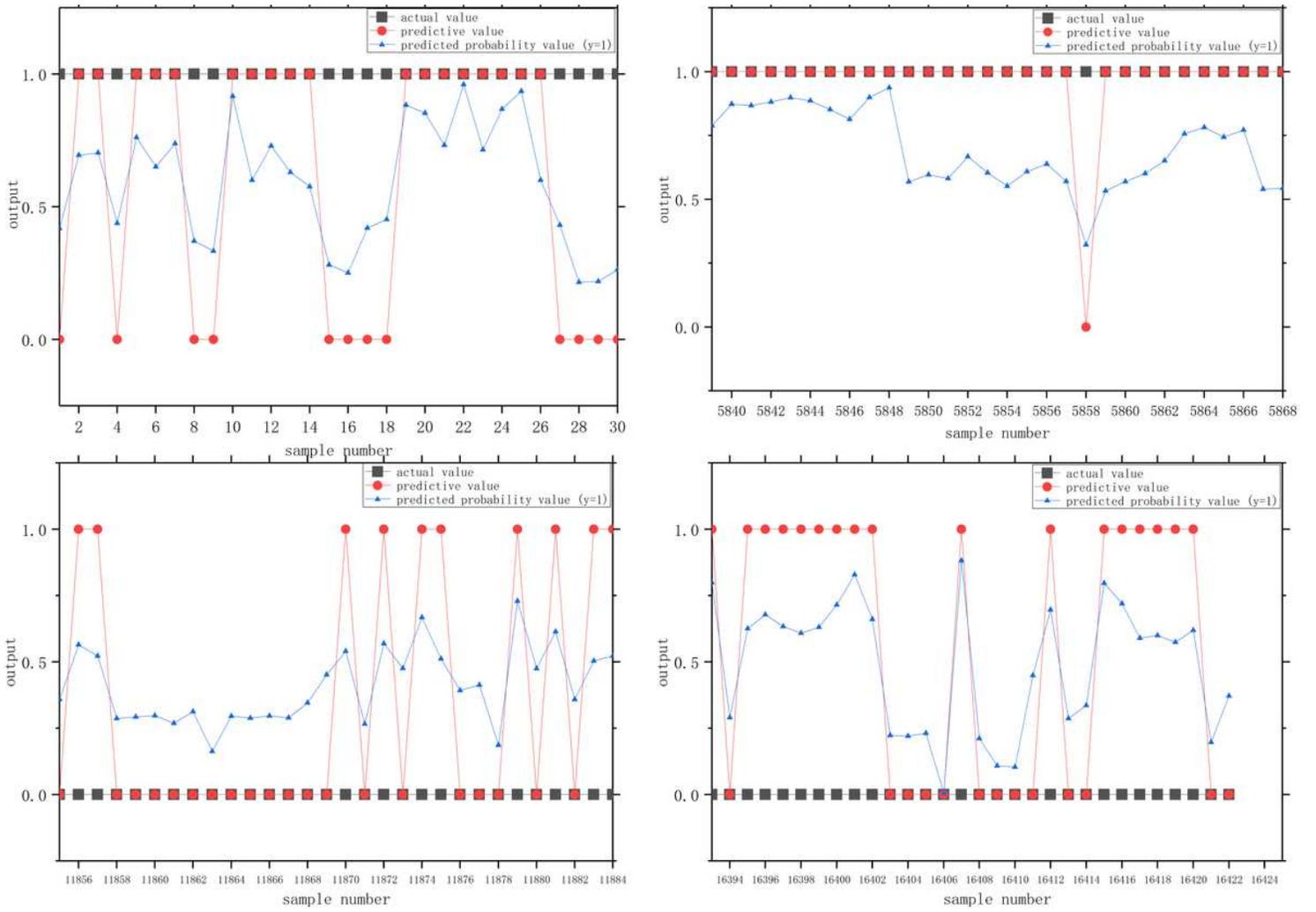


Figure 7

Comparison charts of the predictive and actual values of the RBFNN (part of the sample)

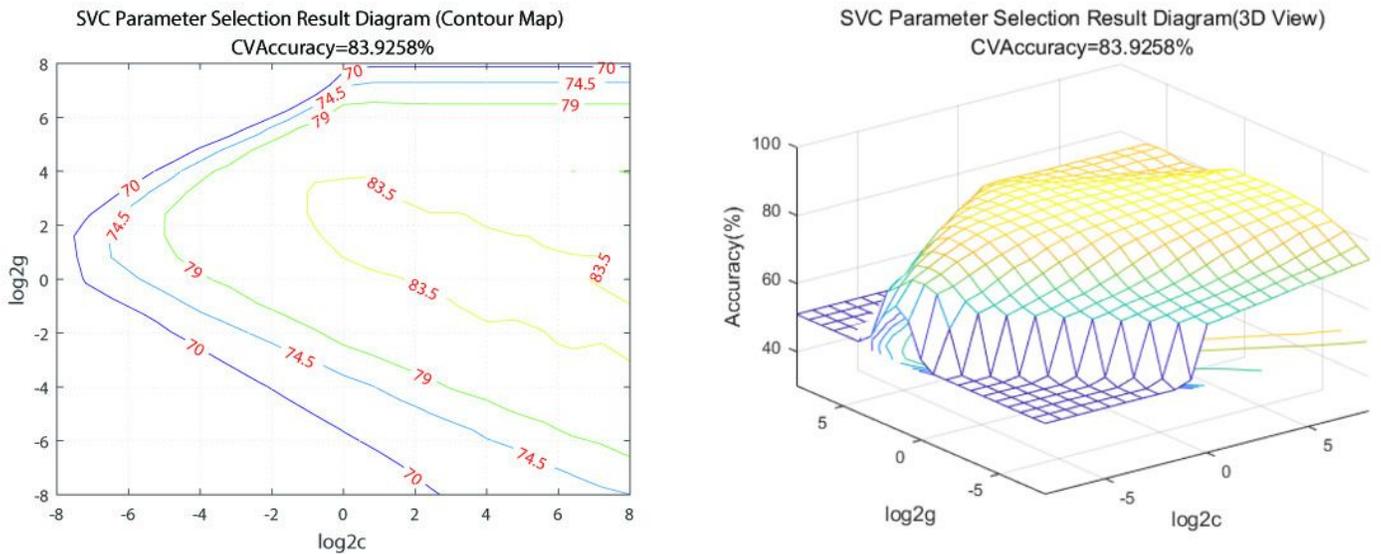


Figure 8

SVC parameter selection result: (a) contour map (b) 3D view

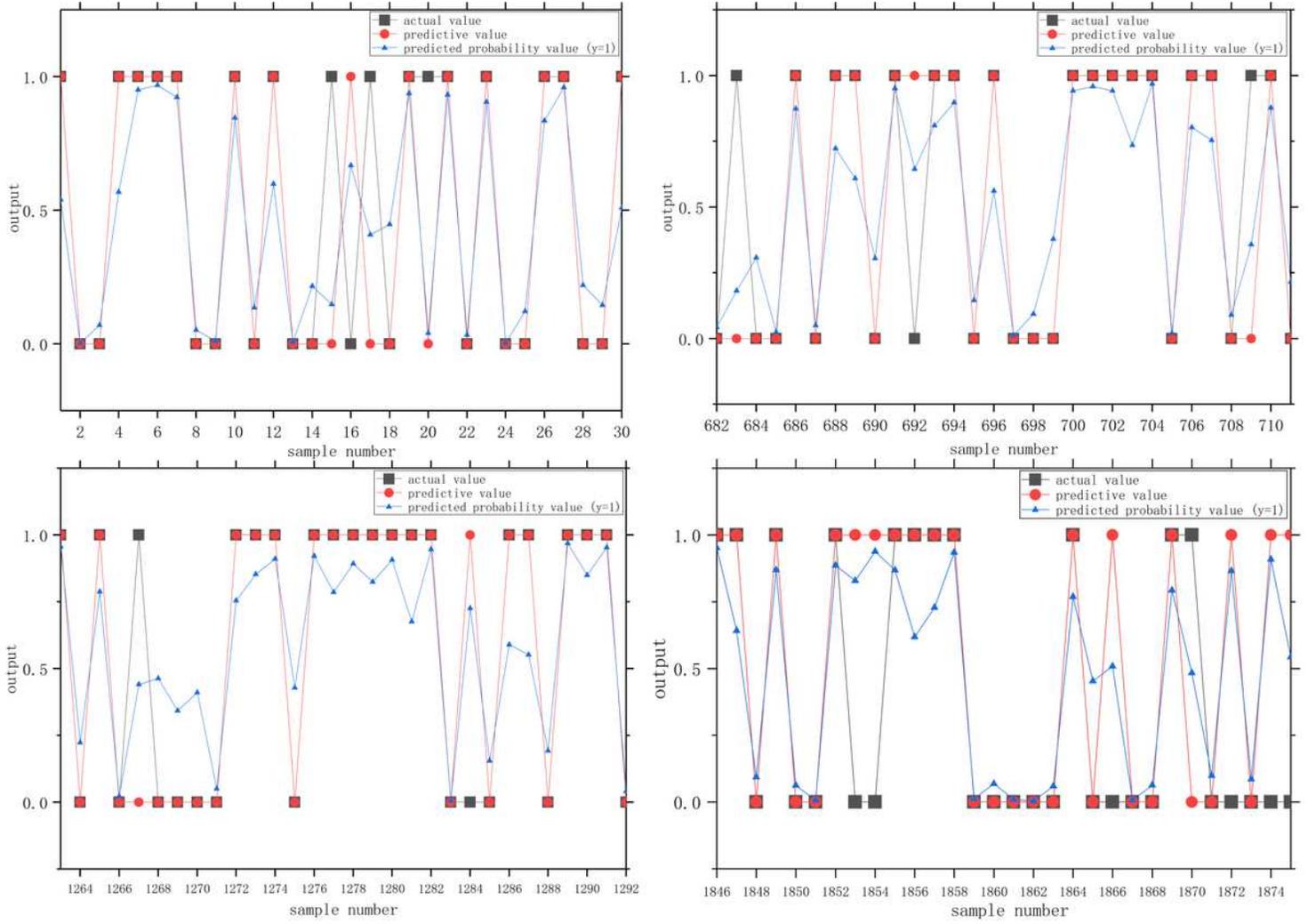


Figure 9

Comparison charts of the predictive and actual values of the SVM (part of the sample)

The number of trees and the accuracy of training and testing in cross-validation

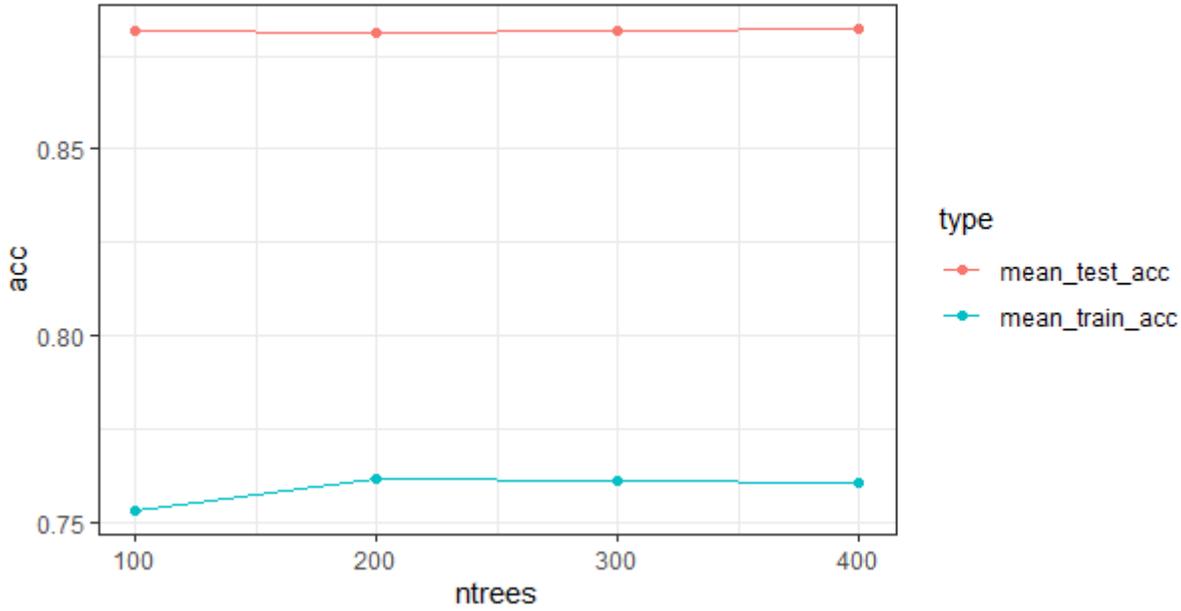


Figure 10

The number of trees and the accuracy of training and testing in cross-validation

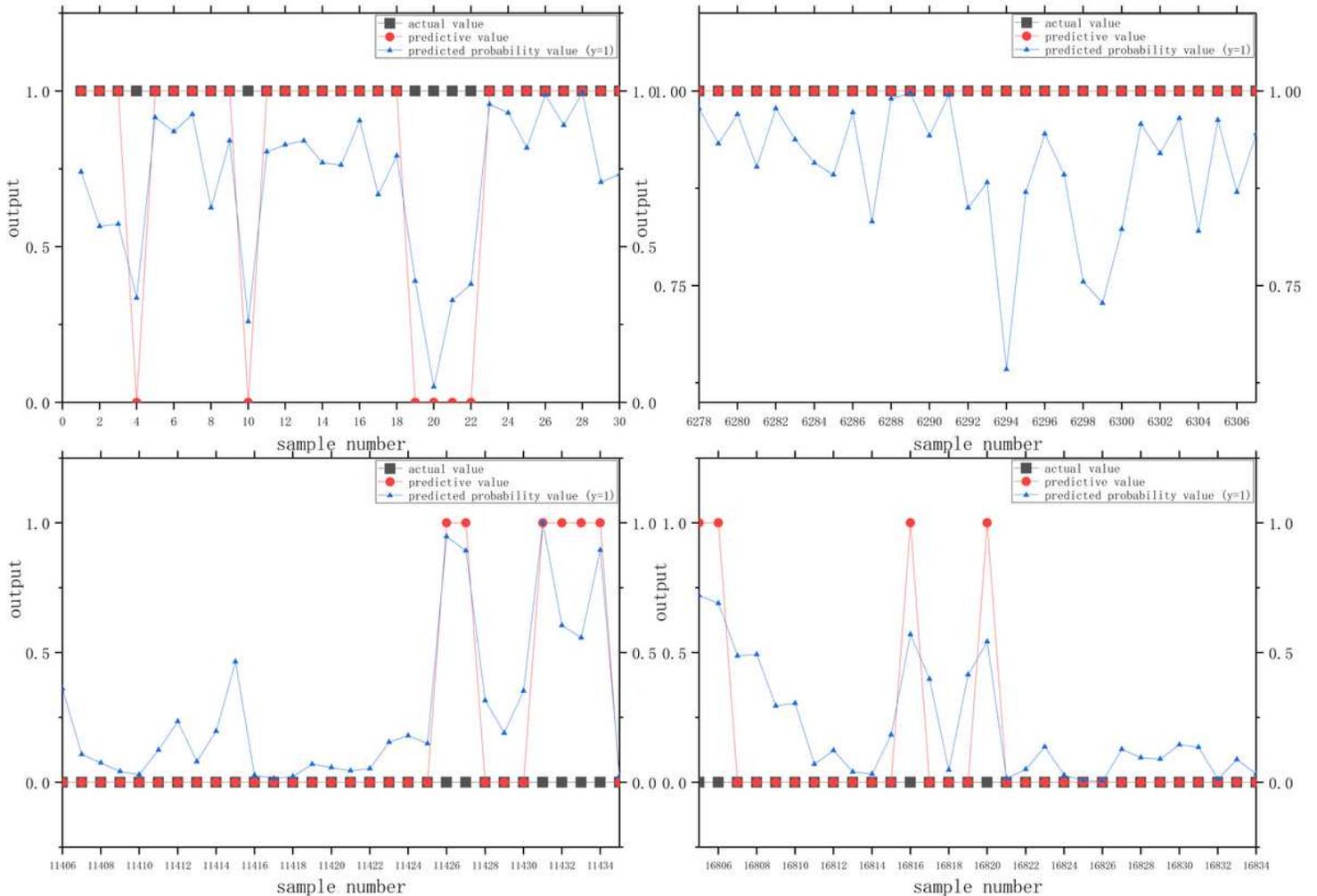


Figure 11

Comparison charts of the predictive and actual values of the RF (part of the sample)

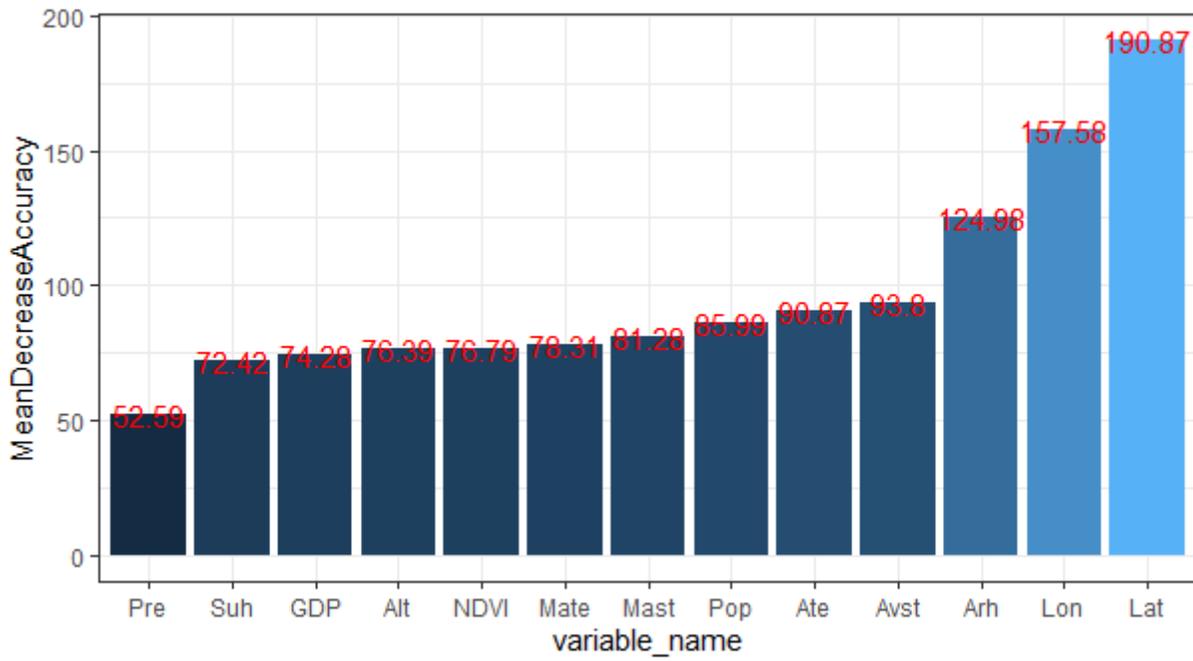


Figure 12

Mean decrease accuracy of 13 variables.

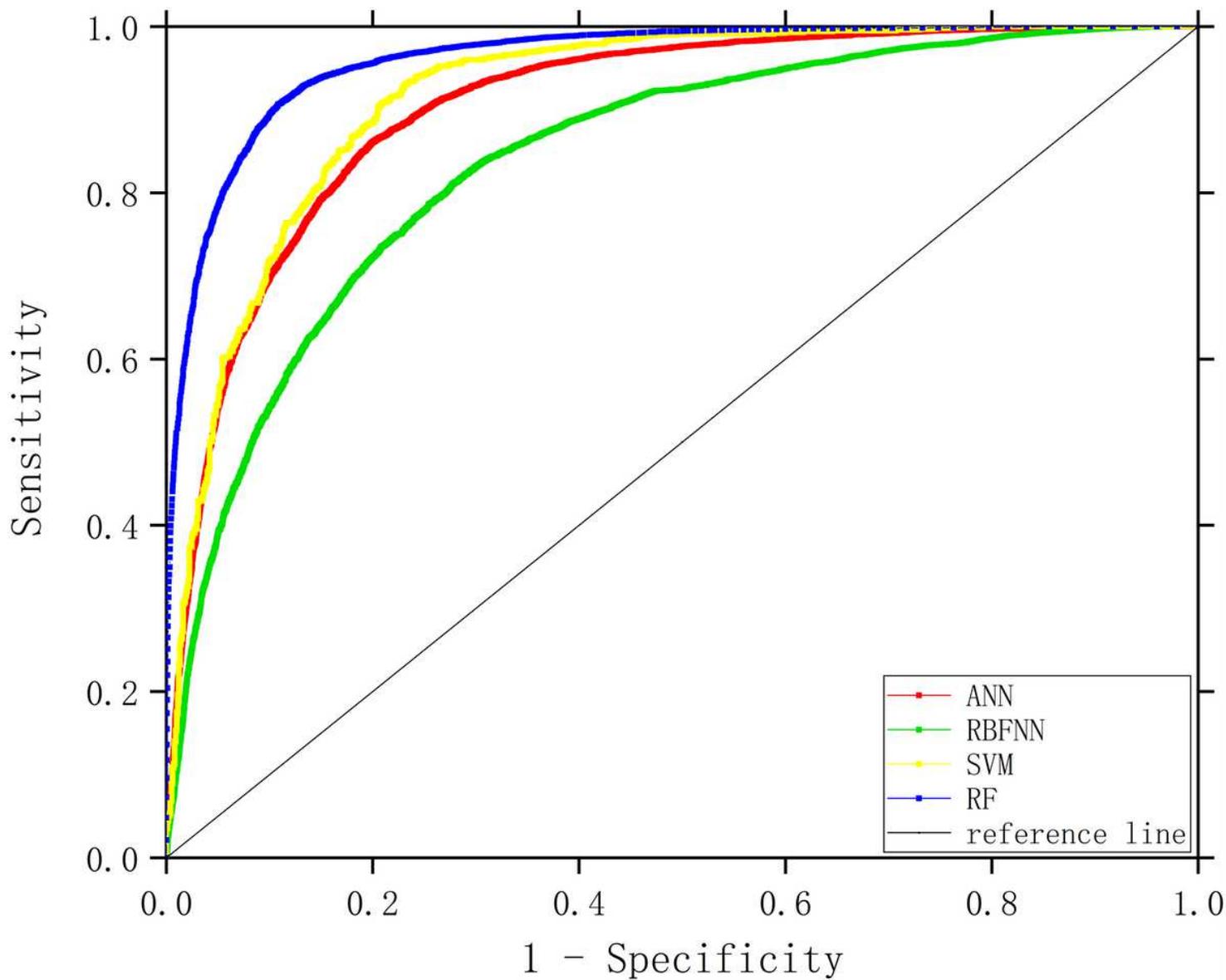


Figure 13

ROC curves of the four models.

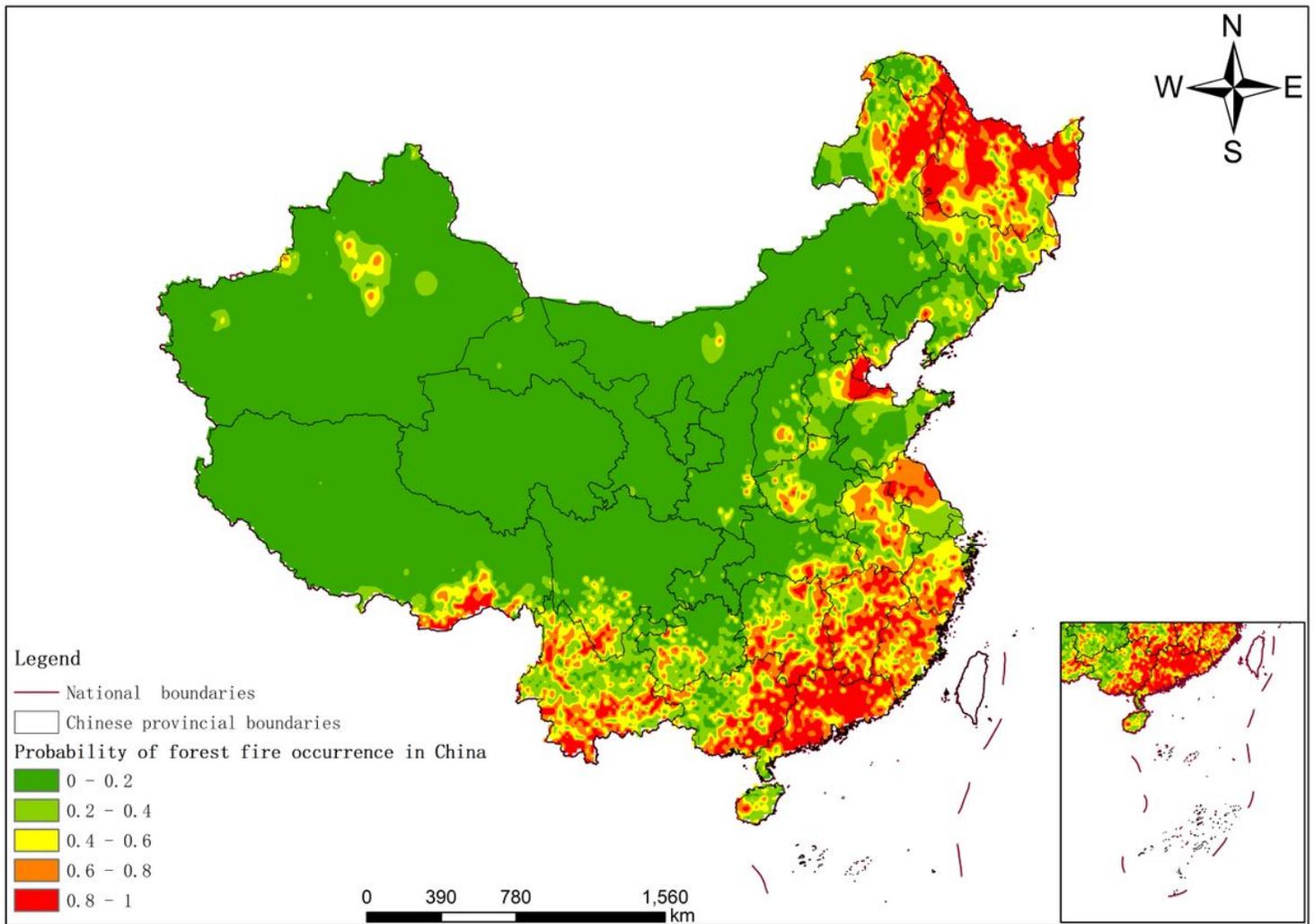


Figure 14

Forest fire probability map for China based on the RF model. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

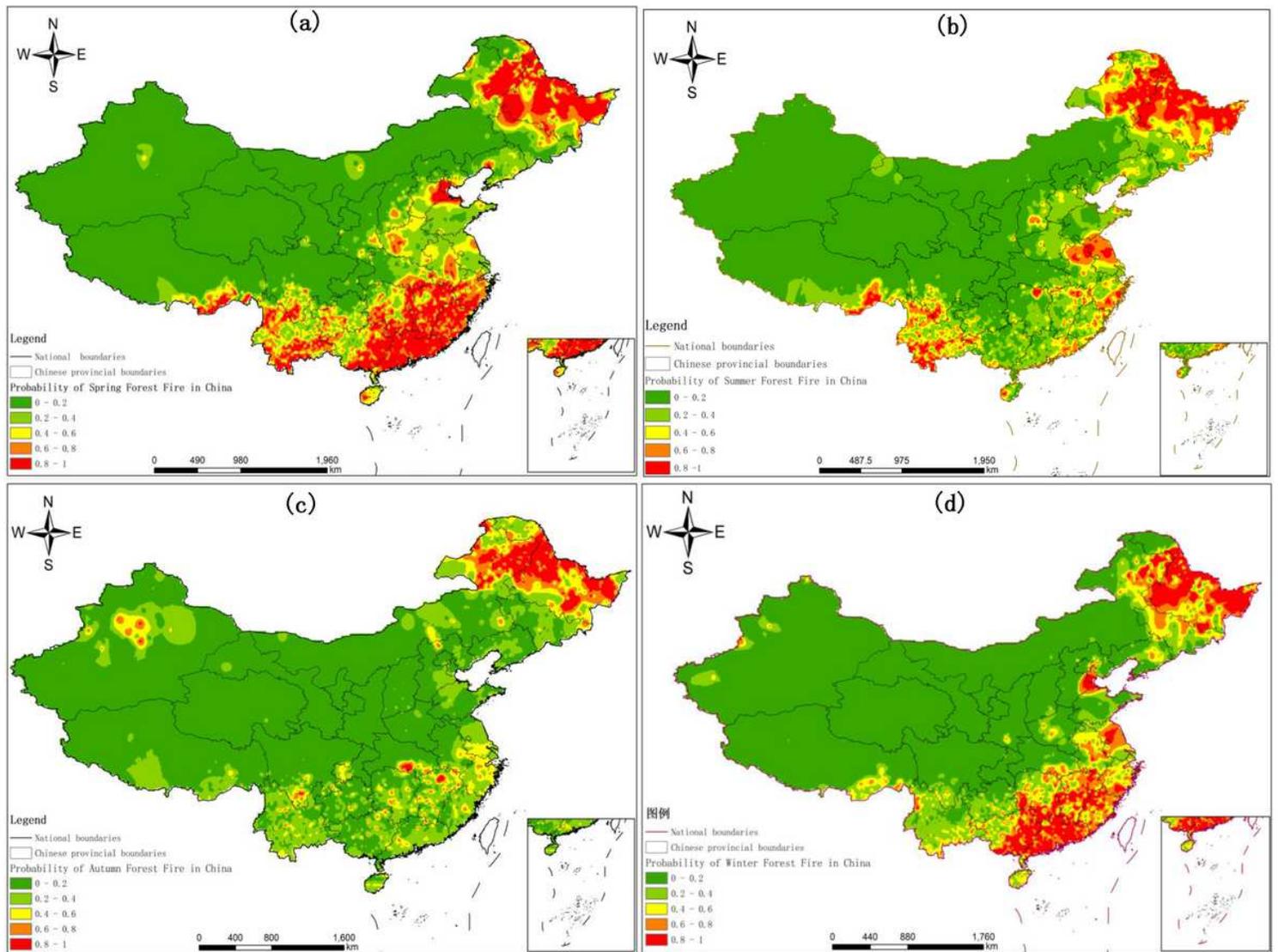


Figure 15

Seasonal forest fire probability map for China based on the RF model: (a) spring (January, February, and March); (b) summer (April, May, and June); (c) autumn (July, August, and September); and (d) winter (October, November, and December). Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.