

Wine recommendation algorithm based on partitioning and stacking integration strategy for Chinese wine consumers

Weisong Mu, Yumeng Feng, Haojie Shu, Bo Wang, Dong Tian*

College of Information and Electrical Engineering, China Agricultural University, Beijing 100083, P.R. China

*Corresponding Author: Dong Tian, College of Information and Electrical Engineering, China Agricultural University, Beijing 100083, P.R. China. Email: td_tiandong@cau.edu.cn

Received: 6 April 2022; Accepted: 10 June 2022; Published: 28 June 2022

© 2022 Codon Publications

OPEN ACCESS 

PAPER

Abstract

This study tries to propose a wine recommendation algorithm based on partitioning and Stacking Integration Strategy for Chinese wine consumers. The approaches follow the idea of partitioning, decomposing traditional recommendation task into several subtasks according to wine attributes, using neural network, support vector machine (SVM), decision tree, random forest, optimized random forest, Adaboost and XGBoost as recommendation models. Then, based on Stacking integration method, five models are screened out for each recommendation index as the base classifier, and the decision tree or logistic regression model is selected as the meta-learner to construct a two-layer Stacking integration framework. Finally, the optimal recommendation algorithm be built for recommendation subtasks according to the prediction accuracy. The result showed that the Stacking integrated recommendation model was suitable for the recommendation of eight attributes including colour, sweetness, foamability, mouthfeel, aroma type, year, packaging and brand, while SVM model was suitable to recommend aroma concentration and price, and the XGboost model was most appropriate for origin. This study would subserve consumers to choose the wine more easily and conveniently and provide support for wine companies to improve customer satisfaction with consumer services. The study expands the approach of concerning research and proposes a specific multi-model recommendation strategy based on artificial intelligence models to recommend multiattribute commodities.

Keywords: Chinese consumers; machine learning; preference for wine attributes; Recommendation algorithm; Stacking integration

Introduction

The development and application of big data analysis technology play an important role in reducing operating costs for businesses and accurately recommending products for customers. Due to the rich attributes and various categories of wine, and indifference of wine knowledge among Chinese consumers, there will be problems such as indecision, time consumption and high cost of trial and error when consumers purchase wine. Therefore, accurate identification of customers' wine drinking preferences

and making personalised recommendations are of great significance to improving customer satisfaction and business performance.

Recommendation algorithm can help users quickly and accurately screen out the corresponding business information from massive data without knowing their own preferences (Liu *et al.*, 2017). Traditionally, similarity measurement and its variants are usually used to make recommendations. With the emergence of machine learning methods in the field of artificial intelligence,

scholars have introduced machine learning technology into recommendation algorithms (Zhang and Lei, 2021), which mainly focuses on using supervised learning algorithms to predict user preferences in different fields (Portugal *et al.*, 2018). A study in 2018 proposed an innovative association classification method, which mined demand satisfaction rules according to the feedback behaviour of users on the recommended products by the system, with the aim of making recommendations based on consumer initiative decision. However, there is a problem of cold start of users, and this method is not applicable when there is no interactive data of users on wine purchase (Yin *et al.*, 2018). Meng and Xiong (2021) used Latent Dirichlet Allocation topic model to divide all user–doctor history consultation text in online medical community into different themes. Then, relevant therapists for users were recommended by combining with the implied disease information in the consultation text and collaborative filtering framework. Compared with the traditional recommendation algorithm, the overall accuracy has improved, whereas the applicable data type is only unstructured text data. Most of the above-mentioned studies used a single machine learning model and optimized it to find out a better prediction result. Although the prediction method of single model is relatively mature, the generalisation ability still needs to be improved. In addition, random factors will affect the model effect and lead to low prediction accuracy (Xie *et al.*, 2020).

In order to effectively reduce or offset the influence of random factors and improve the prediction accuracy and credibility of the prediction model in a single model, ensemble learning has become one of the hot topics in machine learning (Hu *et al.*, 2021). Ensemble methods usually use multiple weak classifiers to form a strong classifier algorithm to improve the classification accuracy. Ensemble learning algorithms mainly include Bagging (He *et al.*, 2019), Boosting and Stacking algorithms, among which Bagging and Boosting algorithms can only integrate single learners of the same type to reduce variance and deviation (Wang *et al.*, 2019). Different from the previous ideas, Stacking can integrate many different types of learners, which can compensate the weakness of single algorithm, enhance the generalisation ability, reduce the risk of overfitting, and improve accuracy. Xie *et al.* (2020) put forward a Stacking framework, which can realize the automatic classification of six types of *Anoectochilus roxburghii* leaves. Considering that the research object is image data with poor interpretability, this method is not suitable for many mixed data sets. Li and Zhai (2019) used Adaboost and random forest as base classifiers of Stacking algorithm to predict and analyse the turnover factors of enterprise staff, aiming at strengthening the management and control of top managers on staffs. Tao *et al.* (2019) used the stacking algorithm to classify the continued spectral series of

rape vegetation, and then embedded the algorithm into unmanned aerial vehicle to distinguish different growing rapes accurately. Based on the interactive behaviour data between users and retail products, Zhang and Lei (2021) established a Stacking algorithm which combines LightGBM, XGboost and random forest to predict the possibility of buying by users in the future and the specific purchase time. However, new users who have not generated behaviour are not considered, it could not solve the cold start problem of new users. Based on the above analysis, in most cases, the combined model has better prediction performance than the single model.

Numerous studies have been conducted on the issue of consumer preference on wine attributes in recent years, which focus on attribute preference influencing factors (Gustafson *et al.*, 2016; Mehta and Bhanja, 2018; Szolnoki and Hauck, 2020) and the actual impact on consumers (Areta *et al.*, 2017; Lee and Lee, 2008; Li *et al.*, 2022). Most of these studies use simple statistical methods or single machine learning algorithms. For example, by using a simple best–worst method, Stanco *et al.* (2020) concluded that traditional attributes such as ‘geographical indications’ and ‘grape variety,’ influence consumer purchase behaviour more than untraditional attributes such as ‘alcohol-free wine’ and ‘vegan wine.’ Chu *et al.* (2020) built a predictive model for Chinese wine consumers’ sensory preferences based on multivariate disorder logistic regression method. On the other hand, using simple logistic regression algorithm is one-sided in analysing the influence of consumers’ personal characteristics on wine selection. Although there has been a great deal of research into consumer preferences on wine, recommendation algorithms and personalised recommendation strategies for wine are less researched and still need further development.

In summary, considering the rich attributes and factors on wine recommendation, and the immature research on personalised recommendation strategy, this study first decomposes the recommendation task of wine product into several recommendation subtasks according to wine attributes with the idea of divide and conquer. Based on the discrete and classified data on wine consumers’ preferences, decision tree, random forest, optimized random forest, neural network, support vector machine (SVM), Adaboost and XGBoost algorithms are used to train the data of each recommendation subtask. The node partition of random forest has a significant effect on the classification performance of multi-classification attributes; therefore, the node partition of random forest algorithm is optimized by linear combination of information gain rate and Gini coefficient, which is used as a recommendation model. Further, the models with the top five accuracies are taken for Stacking integration, and the best recommendation models of each subtask were screened

out according to precision. Finally, these models are used to determine which wine attributes can be recommended to consumers, to establish a specific strategy for wine recommendation.

Materials

Acquisition and preprocessing of data

Preliminary selection of wine recommendation subtasks and independent variables

According to the sensory system of Chinese wine and the research literature of wine field, this study used the inherent wine attributes as recommended index, including colour, sweetness, foamability, mouthfeel, price, concentration and type of aroma, year, origin, packaging and brand of wine. The data types were all discrete. For the convenience of subsequent calculation, the dependent variables were coded by natural numbers according to the specific labels in the recommended index. The inherent attributes of wine applied in this study can be concluded in Figure 1.

Scholars have shown that gender, age, education, occupation, nationality, region and other factors (Capitello *et al.*, 2019; Pagan *et al.*, 2021; Rodríguez-Donate *et al.*, 2021) affect consumers' purchasing behaviour of wine. The research by Cai *et al.* (2015) showed that people with different eating habits have different degrees of aroma

recognition for wine, and this affect consumers' purchase to some extent. For wine, consumer's characteristics is a vital subdivision variable, and these variables interact with each other to influence consumers' purchase of wine (Cai *et al.*, 2015). Based on the synthesis of existing research literature, the factors influencing consumer's purchase behaviour, that is, independent variables proposed in this study, can be divided into the following: personal characteristics, geographical location, psychological and dietary. Among them, psychological factors which affect consumers' wine purchase include the degree of influence of others, advertising, purchase channels, magazines and promotions. The labels can be divided into three categories: high, medium and low influence. The data types of independent variables are all discrete. For the convenience of subsequent calculation, the labels of independent variables were coded by natural numbers, which have been summarised in Table 1.

The data involved in this study originate from the survey of National Grape Industry Technology System, based on a questionnaire survey for Chinese wine drinkers, with a total of 3421 samples collected from June to October 2020. The study was conducted by interviewing a convenient sample of general wine consumers over the age of 18 (regulations indicate that alcohol operators are not allowed to sell alcohol to minors, and 'minors' are defined as natural persons under the age of 18). The participants were invited through different information and social media platforms (web links, WeChat, QQ app). The

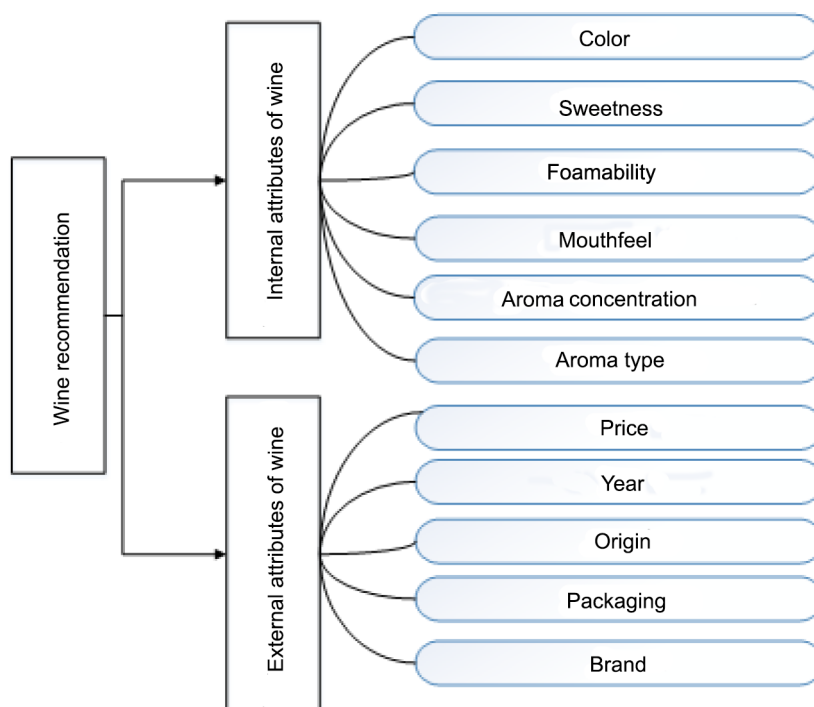


Figure 1. Recommended index of wine.

Table 1. Independent variables of wine recommendation.

Independent variables	Attribute value of independent variable	Attribute value after coding by natural number
Gender	Male, female	(1, 2)
Age	18–25, 26–35, 36–45, 46–55, above 55	(1, 2, 3, 4, 5)
Education level	Below senior high school, High school or technical secondary school, University or Junior college, Postgraduate or more	(1, 2, 3, 4)
Occupation	Students, farmers, freelancers, retirees, employees of state-owned enterprises, employees of private enterprises, party and government organs and institutions, education and scientific research personnel, others	(1, 2, 3, 4, 5, 6, 7, 8, 9)
Monthly disposable income	Under 5000, 5000–10,000, 10,000–15,000, above 15,000	(1, 2, 3, 4)
Marital status	Married, unmarried	(1, 2)
Number of families	3 or less, 4–7, more than 7	(1, 2, 3)
Nation	Han, ethnic minorities	(1, 2)
Address	East China, Northeast China, North China, Central China, South China, Northwest China, Southwest China	(1, 2, 3, 4, 5, 6, 7)
Cognitive level of wine	High, medium, low	(1, 2, 3)
Taste preference	Light taste, heavy taste	(1, 2)
Eating preferences	Sour, sweet, bitter, spicy	(1, 2, 3, 4)
Meat type	White meat, red meat	(1, 2)
Meat preference	Fat meat, lean meat	(1, 2)
Influence of others	High, medium, low	(1, 2, 3)
Influence of magazines	High, medium, low	(1, 2, 3)
Influence of advertisements	High, medium, low	(1, 2, 3)
Influence of promotions	High, medium, low	(1, 2, 3)
Influence of purchase channels	Supermarkets, stores, online stores and chateau	(1, 2, 3, 4)
Data collection		

questionnaires were anonymous and distributed in proportion to the population of each province in China. The sample size obtained from the survey was able to cover the main provincial divisions of China, and at the same time, the sample size supported the operation of machine learning algorithms, and all information collected was used for the purpose of the study. Data preprocessing mainly included optimization of outlier and processing of missing values. Furthermore, after data cleaning, the number of valid questionnaires was 3398.

Independent variables of wine recommendation subtasks

Correlation analysis and significance analysis of variables

1. Mutual information method

Mutual information is mainly used to judge the information contributed by the appearance of one variable to the appearance of another variable (Pascoal *et al.*, 2017). The result of mutual information method is called mutual information value, abbreviated as MI. The value of MI can reflect whether discrete variables are related or not,

and further, it can be used to realize data fusion. The formula is shown in Equation (1):

$$I(x,y) = \sum_{y \in Y} \sum_{x \in X} P(x,y) \log \left(\frac{P(x,y)}{p(x)p(y)} \right) \quad (1)$$

2. Gini coefficient method

Gini coefficient algorithm represents the interaction between the independent variables and wine attributes and judges the significance by the purity change after splitting the attributes (Hao *et al.*, 2021), and finally the value of the importance of each independent variable can be obtained. The formula is shown in Equation (2):

$$\text{Gini}(s) = 1 - \sum_{i=1}^m p_i^2 \quad (2)$$

Eventually independent variables used for recommendation subtasks

In this study, there are 19 independent variables and 11 recommended indicators involved in correlation analysis.

The threshold of MI is 0.01. If the absolute value of MI of two variables is greater than or equal to the threshold, it is considered that the two variables are related, and if the value is less than the threshold, it is judged that the two variables are not related. To further reduce the number of independent variables, Gini coefficient and information gain are used to analyse the importance of independent variables of each recommended index, and finally the top 10 crucial independent variables of each recommended index are obtained. The wine recommendation indicators and the corresponding top 10 independent variables are mentioned in Table 2.

Methodology

Prediction algorithms for wine consumers' sensory attribute preferences

The variables involved in wine recommendation index are all discrete data with obvious labels. Because supervised classification algorithm can effectively solve such problems, single hidden layer BP neural network algorithm, SVM algorithm, decision tree algorithm, random forest algorithm, Adaboost algorithm and Xgboost algorithm are used in the selection of machine algorithms as recommendation models. Table 3 introduces the pros and cons of the selected algorithms for this study.

Random forest is a widely used and powerful classification algorithm currently. Its construction process is roughly as follows: firstly, build multiple decision trees, secondly, use multiple decision trees to train samples, and finally, predict the category of samples in a way that the minority obeys the majority (Janitza *et al.*, 2016). Among them, the construction of decision tree is the vital factor that determines the category of tested samples, and so how to accurately split the nodes is a crucial factor to optimize the random forest. Each recommendation subtask of this data has diverse attributes when building a tree, and the number of classification categories is also varied, so the internal confusion of samples in each subtasks is different. The tree-based models in random forest algorithm involves an optimization problem as to how to split the internal node in chaotic samples and make the results accord with the actual situation. It is necessary to optimize the node-splitting mode in random forest algorithm because the traditional node-splitting standard may cause inaccurate feature division. The CART tree using the Gini coefficient is insufficient for the trend of node division in multi-classified datasets, yet it has higher accuracy of node division for purer datasets. ID3 tree using information gain is easier to distinguish the value of features compared with Gini coefficient when dealing with chaotic datasets, still information gain is partial to the attributes with more values, and as a result, multi-valued attributes are regarded as

Table 2. Recommended indicators of wine and corresponding independent variables.

Recommended subtasks	Eventually independent variables used for modeling (Importance from high to low)
Colour	Meat type, occupation, monthly disposable income, influence of others, age, influence of promotions, advertising influence, education level, influence of magazines, eating preferences
Sweetness	Meat type, monthly disposable income, occupation, influence of others, age, education level, advertising influence, influence of promotions, eating preferences, meat preferences
Foamability	Monthly disposable income, occupation, influence of others, advertising influence, promotion influence, magazine influence, age, education level, meat type, meat preferences
Mouthfeel	Occupation, monthly disposable income, influence of others, advertising influence, age, promotion influence, magazine influence, education level, meat type, eating preferences
Aroma concentration	Occupation, monthly disposable income, advertising influence, influence of others, age, promotion influence, education level, meat preferences, eating preferences, taste preference
Aroma type	Occupation, monthly disposable income, promotion influence, influence of others, advertising influence, age, meat type, eating preferences, education level, meat preferences
Price	Education level, meat type, gender, promotion influence, magazine influence, marital status, advertising influence, age, occupation, monthly disposable income
Year	Monthly disposable income, meat preferences, age, magazine influence, influence of others, advertising influence, eating preferences, gender, promotion influence, marital status
Origin	Influence of others, advertising influence, meat type, monthly disposable income, promotion influence, meat preferences, marital status, eating preferences, age, education level
Packaging	Gender, marital status, influence of others, occupation, advertising influence, education level, magazine influence, meat preferences, eating preferences, meat type
Brand	Magazine influence, monthly disposable income, meat type, meat preferences, eating preferences, education level, age, marital status, advertising influence, gender

Table 3. Strengths and weaknesses of different types of algorithms.

Algorithm	Merit	Defect
BP neural network algorithm	Strong learning ability, easy training, convenient and fast.	Unstable network structure, low reliability and slow convergence (Diez <i>et al.</i> , 2019)
Support vector machine	The effect of processing discrete data is better; The boundaries are more diverse; Solve the problem of non-linear separable classification.	Poor interpretability, slow calculation and easy overfitting.
Decision tree algorithm	Time complexity is small, efficiency is high, importance of each feature can be obtained (Yang and IEEE, 2019), suitable for small samples.	It is easier to make mistakes for tasks with many categories; The information gain is more inclined to the characteristics of multi-category attributes; Prone to overfitting; Poor generalisation ability.
Random forest algorithm	Fast processing speed; For unbalanced data sets, the error can be balanced; Reduce the possibility of overfitting.	Compared with decision tree, the calculation cost is high; When the number of noisy data is large, the results are easily affected by extreme values (Roy and Larocque, 2012) and easy to overfit. The feature weights are biased towards the features with many categories.
Adaboost	No overfitting phenomenon; The training error rate decreases with the increase of iteration times.	The training process is more inclined to the samples that are difficult to classify, which is easily disturbed by noise (Liao and Zhou, 2012); Relying on weak classifier; long training time.
Xgboost	Strong accuracy; Using regular terms to reduce overfitting; good interpretability.	The complexity of time and space is high.

the optimal partition nodes, and the nodes that really need to be partitioned are ignored. Ulteriorly, the information gain ratio further improves this problem by adding term weight to the information gain of each attribute. Therefore, this study structures the linear combination of information gain ratio and Gini coefficient as a new rule of node splitting, to improve the accuracy of random forest algorithm in classification of wine attributes. The formula of the new rule of node splitting is shown in Equation (3):

$$\Phi(\alpha) = \alpha_1 \text{Gini}_{\text{split}}(S, A) - \alpha_2 \text{GainRatio}(A) \quad (3)$$

Where, $0 \leq 1$, the and cannot assume the minimum or maximum value, which means that the weight coefficients cannot exist in such a pairing as (0,0) or (1,1), and the optimal values of the two weights are determined by the grid search method. When dividing the attributes, the node with the smallest difference value is selected for division.

Construction of Stacking integration algorithm

Stacking is a special and concrete combination strategy, which is a typical representative of ensemble learning algorithm (Montesinos-López *et al.*, 2019). The process of Stacking integration is to train the upper layer model first, and take the corresponding training results as the input of the next layer model, which is equivalent to constructing a two serial superposition classifier for the same data (Ahmadi *et al.*, 2019). The purpose is to improve the classification accuracy of the model.

In this study, decision tree, random forest, optimized random forest, neural network, SVM, Adaboost and Xgboost algorithms modeled each wine recommendation subtask, and the top five models with the highest recommended comprehensive accuracy of each recommendation index were selected as the base models. Further, to avoid the overfitting caused by complicated model, the second meta-model selected decision tree or logistic regression model. The construction process of Stacking integrated algorithm can be concluded in Figure 2.

Multimodel wine recommendation strategy

To establish the most appropriate recommendation strategy of each recommendation subtask, it is necessary to compare the recommended comprehensive accuracy of the models before and after Stacking integration. The final framework of wine recommendation in this study is shown in Figure 3.

The personalised wine recommendation method comprises the following steps:

Step 1: Collect user's preference data for wine, including 19 independent variables and 11 dependent variables.

Step 2: Clean the data, eliminate or optimize the wrong information, and encode the data by natural number after obtaining the effective information.

Step 3: Experiment with decision tree, random forest, optimized random forest, neural network, SVM, Adaboost, Xgboost and Stacking integration model for each recommendation subtask and select the model with

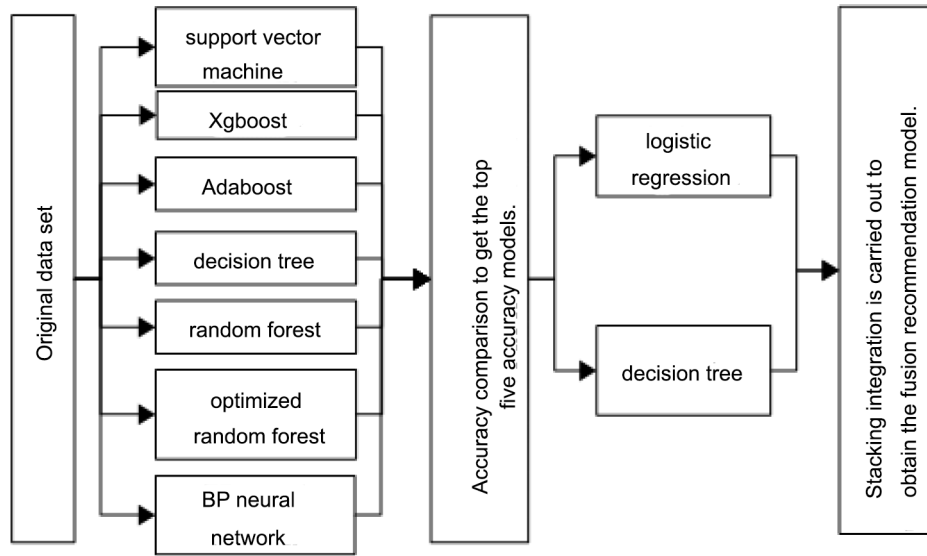


Figure 2. The construction process of Stacking integration model.

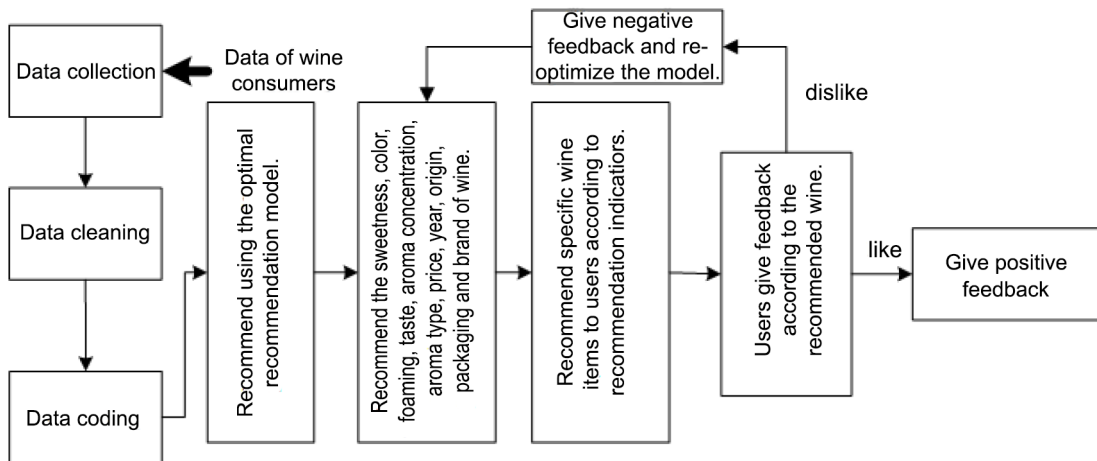


Figure 3. The process of personalised wine recommendation.

the highest accuracy as the optimal recommendation strategy to recommend wine.

Step 4: There is a user interaction module in the recommendation method. If the user likes the recommended wine, they give correct feedback to the recommendation model; if the user doesn't like it, they give negative feedback so that the strategy can continue to learn and further be optimized.

Evaluation index for recommendation algorithm

Accuracy

Accuracy refers to the proportion of the sum of correctly classified samples to the total samples in

the prediction results. The formula is expressed in Equation (4):

$$acc = \frac{TP + TN}{total} \tag{4}$$

m-F1

The F1 score considers both the accuracy and recall of the classification model and can be regarded as a harmonic mean of the accuracy and recall. In this study, we put forward the concept of m-F1, which represents the arithmetic average of F1 scores of each category of the same model and the same recommended subtask.

Table 4. The recommended comprehensive accuracy of each single model.

	Colour	Sweetness	Foamability	Mouthfeel	Aroma type	Brand	Price	Year	Origin	Packaging	Aroma concentration
BP neural network	0.788	0.755	0.801	0.765	0.779	0.756	0.733	0.757	0.773	0.770	0.741
Support vector machine (SVM)	0.802	0.827	0.791	0.717	0.761	0.771	0.762	0.741	0.747	0.792	0.777
Decision tree	0.783	0.797	0.758	0.743	0.677	0.753	0.738	0.728	0.732	0.751	0.737
Random forest	0.817	0.829	0.763	0.756	0.689	0.774	0.748	0.742	0.745	0.763	0.762
Optimized random forest	0.833	0.845	0.792	0.786	0.716	0.798	0.761	0.629	0.745	0.738	0.466
Adaboost	0.825	0.842	0.761	0.752	0.686	0.787	0.742	0.736	0.781	0.780	0.774
XGboost	0.818	0.853	0.788	0.771	0.725	0.745	0.742	0.736	0.778	0.784	0.726

Table 5. The base model and meta-model of Stacking integration for each recommended subtask.

Recommended index	Base model 1	Base model 2	Base model 3	Base model 4	Base model 5	Meta-model
Foamability	BP neural network	Optimized random forest	Support vector machine (SVM)	XGboost	Random forest	Logical regression
Mouthfeel	Optimized random forest	XGboost	BP neural network	Random forest	Adaboost	Decision tree
Aroma concentration	SVM	Adaboost	Random forest	BP neural network	Decision tree	Logical regression
Colour	Optimized random forest	Adaboost	XGboost	Random forest	Support vector machine (SVM)	Decision tree
Sweetness	XGboost	Optimized random forest	Adaboost	Random forest	SVM	Logical regression
Aroma type	BP neural network	SVM	XGboost	Optimized random forest	Random forest	Logical regression
Year	BP neural network	Random forest	SVM	Adaboost	XGboost	Decision tree
Origin	XGboost	BP neural network	SVM	Random forest	Optimized random forest	Logical regression
Packaging	SVM	XGboost	Adaboost	BP neural network	Random forest	Logical regression
Brand	Optimized random forest	Adaboost	Random forest	SVM	BP neural network	Decision tree
Price	SVM	Optimized random forest	Random forest	XGboost	Adaboost	Logical regression

Recommended comprehensive accuracy

In this study, the concept of recommended comprehensive accuracy is proposed, which is expressed by the arithmetic average of m-F1 and accuracy.

Results and Discussion

Construction and verification of single recommendation model

This section needs to complete the training and evaluation of single recommendation model of each wine subtask. The independent variables and recommendation indicators used in building the recommendation model are shown in Section 2.2. The 10-fold cross-validation method is used to divide and train the sample. To better evaluate the effect of the model, this study randomly selected 300 questionnaires from the valid data to test the model. The recommended comprehensive accuracy proposed in this study blends the accuracy and m-F1 scores, so that it only presents the recommended comprehensive accuracy of each single model for each wine recommended subtask. The results are shown in Table 4.

The recommendation comprehensive accuracy shows that the optimized random forest is better than the random forest in colour, sweetness, foaming, taste, aroma type, price and brand of wine, and the precision has been improved. For the colour, brand and taste of wine, changing the splitting mode of nodes can improve the recommendation accuracy, and the optimized random forest algorithm model has the best recommendation effect. As for sweetness, there is some noise in the sample. Compared with other subtasks, sweetness of wine pays more attention to variance in modeling, while XGboost adopts regularisation constraint to simplify the model. Using XGboost makes

the comprehensive accuracy up to 0.853. For foamability, aroma type, year and origin of wine, the distribution of these data is diverse, and this kind of data need a better fitting algorithm for recommendation. Among these machine learning algorithms, the advantage of neural network is to constantly fit data. So neural network performs best. For price, aroma concentration and packaging, there is almost no noise, and there are some key points similar to the support vector in the data, so SVM is more suitable for the recommendation of these wine attributes.

Screening of base model for Stacking integration algorithm

According to the above results, the top five recommendation models with the highest comprehensive accuracy of each subtask are selected as the base models. To avoid the overfitting caused by the complexity of the two-layer models, the decision tree or logistic regression model is selected as the meta-model to reduce the complexity. The established machine learning algorithm of the two-layer basic model for each wine recommendation subtask is shown in Table 5.

Implementation and verification of Stacking Integrated Algorithm

Firstly, train the single model and the integrated model on the training set. Secondly, randomly select 300 questionnaires from the valid data as the test set and make predictions. By comparing the recommendation comprehensive accuracy of each subtask before and after Stacking integration, the model with the highest value is selected as the optimal recommendation model, and finally an exclusively multimodel recommendation strategy for wine is established. The recommendation

Table 6. The comprehensive accuracy of the optimal model of each recommended index.

Recommended index	Optimal recommendation model	Comprehensive accuracy after Stacking integration	Optimal Comprehensive accuracy before Stacking integration
Sweetness		0.838	0.833
Colour		0.873	0.853
Foamability		0.820	0.801
Mouthfeel	Stacking integration model	0.803	0.786
Aroma type		0.803	0.779
Year		0.771	0.757
Packaging		0.815	0.792
Brand		0.804	0.798
Price	Support vector machine	0.731	0.762
Aroma concentration		0.760	0.777
Origin	XGboost	0.762	0.778

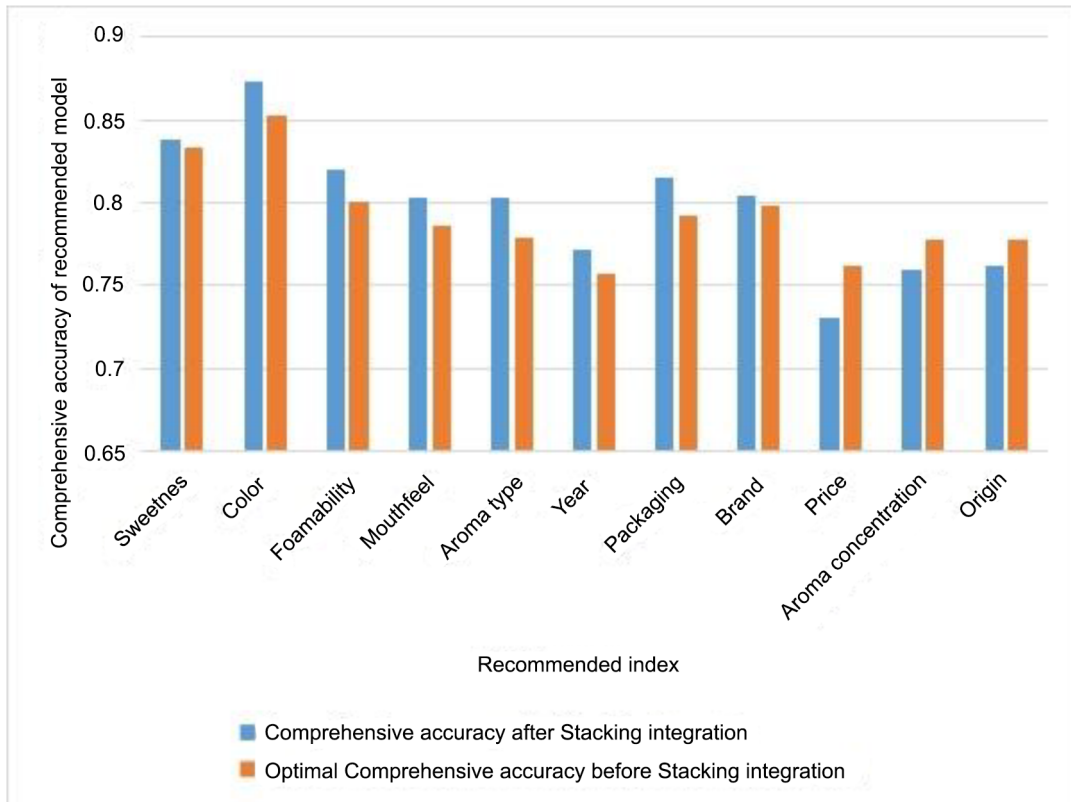


Figure 4. Comprehensive accuracy of recommended model.

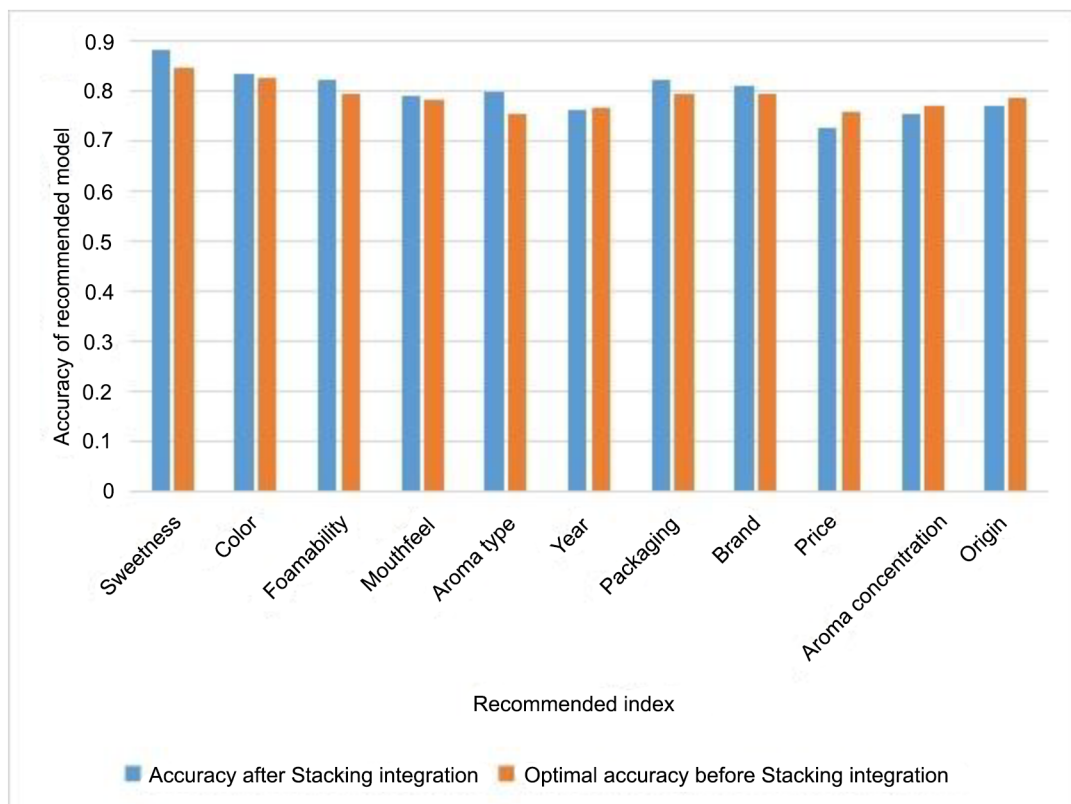


Figure 5. Accuracy of recommended model.

comprehensive accuracy after fusion and the highest value before fusion are summarised as follows.

To present the results more visually, Table 6 is transformed into Figure 4.

As can be seen from Table 6, for the eight indexes of colour, sweetness, foamability, mouthfeel, aroma type, year, packaging and brand, the optimal recommendation model is Stacking integrated algorithm, while the SVM model is most suitable for the recommendation of aroma concentration and price of wine, and XGboost model is the best algorithm for origin. The Stacking integration strategy has better performance than the single model and has significant effect on improving the accuracy of wine recommendation.

The accuracy after fusion and the highest value before fusion are shown in Figure 5.

The literature (Chu *et al.*, 2020) uses a predicting method based on multivariate disorder logistic and shows that this method is only suitable for predicting colour, sweetness, aroma type and taste preference with an accuracy of between 70 and 75%, whereas the wine recommendation algorithm proposed in this study achieved an accuracy of over 79% for these attributes, which is better than the method proposed in the literature.

Conclusions and Future Work

In this study, a multimodel recommendation algorithm strategy is proposed to predict users' preferences about wine and realize personalised recommendation of wine. The independent variables of each wine recommendation subtask were screened by mutual information method and Gini coefficient method, and the top 10 factors with the highest correlation of each recommendation subtasks included personal characteristics of consumers such as age, education level, monthly disposable income and occupation, which indicated that personal characteristics of consumers had a crucial influence on wine purchase.

By combining the information gain rate and Gini coefficient to change the node splitting mode, the random forest algorithm was further optimized, which made the algorithm more accurate in recommending seven indicators such as colour, sweetness, foamability, mouthfeel, aroma type, price and brand of wine, and also showed that the optimized random forest algorithm is more suitable for classification tasks with numerous feature divisions.

To effectively reduce or offset the influence of random factors and improve the prediction accuracy and

credibility in a single model, the Stacking integration was introduced into wine recommendation. The results showed that the integration algorithm could be used to improve the recommendation accuracy of eight indicators, including colour, sweetness, foamability, mouthfeel, aroma type, year, packaging and brand. The average recommendation accuracy is increased by 1.9%, and the highest is increased by 3%. However, it is not suitable for the attributes with large deviation such as aroma concentration. For this kind of tasks, Adaboost and other algorithms are more suitable for recommendation.

Based on the above research, the conclusion of this study is that colour, sweetness, foamability, mouthfeel, aroma type, year and packaging adopted Stacking integrated algorithm, aroma concentration finally adopted Adaboost algorithm, origin finally adopted XGboost algorithm, and wine price finally adopted SVM algorithm. Compared to the existing recommendation methods, the wine recommendation strategy proposed in this study can be more perfect, comprehensive and flexible in recommending wine for users.

Our approach might still be needed to improve, and the future attempts are as follows: Firstly, due to the limited amount of data in this study, the accuracy and the screening results of the optimal model may be varied for a large amount of data, so further exploration is still needed. Secondly, the future work can be expanded to conduct empirical analysis on other multiattribute drinks' recommendation. In addition, one of the typical examples in China's market, the wine multimodel recommendation strategy may also provide an effective method for dealing with preference prediction in other multi-feature goods.

References

- Ahmadi, E., Garcia-Arce, A., Masel, D.T., Reich, E., Puckey, J. and Maff, R. 2019. A metaheuristic-based stacking model for predicting the risk of patient no-show and late cancellation for neurology appointments. *IISE Transactions on Healthcare Systems Engineering*. 9(3): 272–291. <https://doi.org/10.1080/24725579.2019.1649764>
- Arete, A., Bardaji, I. and Iraizoz, B. 2017. Spanish wines in the US market: what attributes do US consumers look for in Spanish wines? *Spanish Journal of Agricultural Research*. 15(4): e0120. <https://doi.org/10.5424/sjar/2017154-10006>
- Cai, J., Pang, J., Hu, K. and Tao, Y. 2015. Chinese familiarity with wine aromas characteristics. *Journal of Food Science and Technology*. 33(04):47–51.
- Capitello, R., Bazzani, C. and Begalli, D. 2019. Consumer personality, attitudes and preferences in out-of-home contexts. *International Journal of Wine Business Research*. 31(1): 48–67. <https://doi.org/10.1108/IJWBR-06-2018-0022>

- Chu, X., Li, Y., Xie, Y., Tian, D. and Mu, W. 2020. Regional difference analyzing and prediction model building for Chinese wine consumers' sensory preference. *British Food Journal*. 122(8): 2587–2602. <https://doi.org/10.1108/BFJ-06-2019-0465>
- Diez, F.J., Navas-Gracia, L.M., Martinez-Rodriguez, A., Correa-Guimaraes, A. and Chico-Santamarta, L. 2019. Modelling of a flat-plate solar collector using artificial neural networks for different working fluid (water) flow rates. *Solar Energy*. 188: 1320–1331. <https://doi.org/10.1016/j.solener.2019.07.022>
- Gustafson, C.R., Lybbert, T.J. and Sumner, D.A. 2016. Consumer sorting and hedonic valuation of wine attributes: exploiting data from a field experiment. *Agricultural Economics*. 47(1): 91–103. <https://doi.org/10.1111/agec.12212>
- Hao, L., Fang, Y. and Zhang, Q. 2021. Rediscussion on statistical measure on Gini coefficient. *Statistics & Decision*. 37(7): 27–32.
- He, Y., Zou, H. and Yu, H. 2019. Recommender system model based on dynamic-weighted bagging matrix factorization. *Journal of Nanjing University (Natural Science)*. 55(04): 644–650.
- Hu, Y., Qu, B., Liang, J., Wang, J. and Wang, Y. 2021. A survey on evolutionary ensemble learning algorithm. *Chinese Journal of Intelligent Science and Technology*. 3(1):18–33.
- Janitza, S., Tutz, G. and Boulesteix, A. 2016. Random forest for ordinal responses: prediction and variable selection. *Computational Statistics and Data Analysis*. 96(C): 57–73.
- Lee, S.J. and Lee, K.G. 2008. Understanding consumer preferences for rice wines using sensory data. *Journal of the Science of Food and Agriculture*. 88(4): 690–698. <https://doi.org/10.1002/jsfa.3137>
- Li, Q. and Zhai, L. 2019. Analysis and research on employee turnover prediction based on stacking algorithm. *Journal of Chongqing Technology and Business University (Natural Science Edition)*. 36(01): 117–123.
- Li, Y., Jia, X., Wang, R., Qi, J., Jin, H., Chu, X. and Mu, W. 2022. A new oversampling method and improved radial basis function classifier for customer consumption behavior prediction. *Expert Systems with Applications*. 199: 116982. <https://doi.org/10.1016/j.eswa.2022.116982>
- Liao, H. and Zhou, D. 2012. Review of AdaBoost and its improvement. *Computer Systems & Applications*. 21: 240–244.
- Liu, H., Guo, M. and Pan, W. 2017. Overview of personalized recommendation systems. *Journal of Changzhou University (Natural Science Edition)*. 29: 51–59.
- Mehta, R. and Bhanja, N. 2018. Consumer preferences for wine attributes in an emerging market. *International Journal of Retail & Distribution Management*. 46(1): 34–48. <https://doi.org/10.1108/IJRDM-04-2017-0073>
- Meng, Q. and Xiong, H. 2021. Doctor recommendation based on online consultation text information. *Information Science*. 39(06): 152–160.
- Montesinos-López, O.A., Montesinos-López, A., Crossa, J., Cuevas J., Montesinos-López, J.C., Gutiérrez, Z.S., Lillemo, M., Philomin, J. and Singh, R. 2019. A Bayesian genomic multi-output regressor stacking model for predicting multi-trait multi-environment plant breeding data. *G3 (Bethesda)*. 9(10): 3381–3393. <https://doi.org/10.1534/g3.119.400336>
- Pagan, K.M., Giraldo, J.D.M.E., Maheshwari, V., de Paula, A.L.D. and de Oliveira, J.H.C. 2021. Evaluating cognitive processing and preferences through brain responses towards country of origin for wines: the role of gender and involvement. *International Journal of Wine Business Research*. 33(4): 481–501. <https://doi.org/10.1108/IJWBR-08-2020-0043>
- Pascoal, C., Oliveira, M.R., Pacheco, A. and Valadas, R. 2017. Theoretical evaluation of feature selection methods based on mutual information. *Neurocomputing*. 226: 168–181.
- Portugal, I., Alencar, P. and Cowan, D. 2018. The use of machine learning algorithms in recommender systems: A systematic review. *Expert Systems with Application*. 97: 205–227. <https://doi.org/10.1016/j.eswa.2017.12.020>
- Rodríguez-Donate, M.C., Romero-Rodríguez, M.E. and Cano-Fernández, V.J. 2021. Wine consumption preferences among generations X and Y: an analysis of variability. *British Food Journal*. 123(11): 3557–3575. <https://doi.org/10.1108/BFJ-12-2020-1156>
- Roy, M. and Larocque, D. 2012. Robustness of random forests for regression. *Journal of Nonparametric Statistics*. 24(4): 93–1006. <https://doi.org/10.1080/10485252.2012.715161>
- Stanco, M., Lerro, M. and Marotta, G. 2020. Consumers' preferences for wine attributes: a best–worst scaling analysis. *Sustainability*. 12(7): 2819. <https://doi.org/10.3390/su12072819>
- Szolnoki, G. and Hauck, K. 2020. Analysis of German wine consumers' preferences for organic and non-organic wines. *British Food Journal*. 122(7): 2077–2087. <https://doi.org/10.1108/BFJ-10-2019-0752>
- Tao, Y., Peng, Y., Jiang, Q., Li, Y., Fang, S. and Gong, Y. 2019. Remote detection of critical growth stages in rapeseed using vegetation spectral and stacking combination method. *Journal of Geomatics*. 44(5): 20–23.
- Wang, Y., Liao, X. and Li, S. B. 2019. Resealed boosting in classification. *IEEE Transactions on Neural Networks and Learning Systems*. 30(9): 2598–2610.
- Xie, W., Chai, Q., Gan, Y., Chen, S., Zhang, X. and Wang, W. 2020. Strains classification of *Anoectochilus roxburghii* using multi-feature extraction and Stacking ensemble learning. *Transactions of the Chinese Society of Agricultural Engineering*. 36: 203–210.
- Yang, F.J. & IEEE. 2019. An extended idea about decision trees. 2019 6TH INTERNATIONAL CONFERENCE ON COMPUTATIONAL SCIENCE AND COMPUTATIONAL INTELLIGENCE (CSCI 2019). 6th Annual Conference on Computational Science and Computational Intelligence (CSCI 2019). 349–354.
- Yin, C., Guo, Y., Yang, J. and Ren, X. 2018. A new recommendation system on the basis of consumer initiative decision based on an associative classification approach. *Industrial Management & Data Systems*. 118(1):188–203. <https://doi.org/10.1108/IMDS-02-2017-0057>
- Zhang, J. and Lei, J. 2021. Stacking Fusion Model for customer purchase behavior prediction. *Shanghai Management Science*. 43: 12–19.