

Short-Term Prediction of Ammonia Levels In Geese Houses Based On Combined Feature Selector and Random Forest

Jiande Huang

Zhongkai University of Agriculture and Engineering

Jianjun Guo

Zhongkai University of Agriculture and Engineering

Huilin Wu

National S&T Innovation Center for Modern Agricultural Industry

Xinglong Zhang

National S&T Innovation Center for Modern Agricultural Industry

Shuangyin Liu

Zhongkai University of Agriculture and Engineering

Shahbaz Gul Hassan (✉ mhasan387@zkhu.edu.cn)

Zhongkai University of Agriculture and Engineering

Article

Keywords:

Posted Date: March 21st, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1449095/v1>

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

[Read Full License](#)

1 Short-Term Prediction of Ammonia Levels In Geese Houses Based 2 On Combined Feature Selector and Random Forest

3 Jiande Huang ^{1,2,3,4,5, †}, Jianjun Guo ^{1,2,3,4,5, †}, Huilin Wu ⁶, Xinglong Zhang ⁶, Shuangyin Liu ^{1,2,3,4,5, *} and
4 Shahbaz Gul Hassan ^{1,2,3,4,5, *}

5 ¹College of Information Science and Technology, Zhongkai University of Agriculture and Engineering, Guangzhou 510225, China

6 ²Smart Agriculture Engineering Technology Research Center of Guangdong Higher Education Institutes, Zhongkai University of
7 Agriculture and Engineering, Guangzhou 510225, China

8 ³Guangzhou Key Laboratory of Agricultural Products Quality & Safety Traceability Information Technology, Zhongkai University of
9 Agriculture and Engineering, Guangzhou 510225, China

10 ⁴Academy of Intelligent Agricultural Engineering Innovations, Zhongkai University of Agriculture and Engineering, Guangzhou 510225,
11 China

12 ⁵Guangdong Provincial Agricultural Products Safety Big Data Engineering Technology Research Center, Zhongkai University of
13 Agriculture and Engineering, Guangzhou 510225, China

14 ⁶National S&T Innovation Center for Modern Agricultural Industry, Guangzhou 510520, China

15 [†]Contributed equally to this work.

16 *Correspondence: Shuangyin Liu (hdlsyxlq@126.com); Shahbaz Gul Hassan (mhasan387@zku.edu.cn)

17 Abstract

18 Ammonia concentration (NH_3) is a dominant source of environmental pollution in geese housing and
19 profoundly affects the healthy growth of geese. Accurately forecasting NH_3 and analyzing its change
20 trends in geese houses is crucial for the survival of geese. To improve the prediction accuracy of NH_3 ,
21 we propose a novel forecasting model by the combination of feature selector (CFS) and random
22 forest (RF). The developed model integrated two modules. First, combining mutual information (MI)
23 and relief-F, we propose CFS quantify the importance values of each feature and eliminate the low-
24 relation or unrelated features. Second, we built a random forest model and used the K-fold cross-
25 validation grid search algorithm (CVGS) to obtain the RF hyper-parameters to predict NH_3 . The
26 simulation results show that the prediction accuracy was improved when feature selection after
27 quantification based on the CFS was used. The mean square error (MSE), root mean square error
28 (RMSE) and mean absolute percent error (MAPE) for the proposed model were 0.5072, 0.6583, 2.88%,
29 respectively. The NH_3 prediction model based on CFS-GS-RF exhibited best prediction accuracy, and
30 generalization performance compared with other parallel forecasting models and is a suitable and
31 useful tool for predicting NH_3 in geese houses.

32 33 Introduction

34 China is the largest geese producer in the world. According to the data published by the China
35 Statistics Bureau in 2019 National Economic and Social Development Statistical Bulletin, the output
36 of poultry meat and poultry eggs in China were 22.39 million tons and 33.09 million tons,

37 respectively. The total output value of waterfowl exceeds 160 billion yuan ¹. Due to the increasing
38 consumption of poultry products and an increase in exports, more poultry coops will be built.
39 Historically, outbreaks of poultry disease and even mass deaths are almost inevitable due to rapid
40 expansion of breeding scale, lack of scientific management and environmental degradation of the
41 poultry housing ². Presently, NH₃ produced in poultry houses poses the most significant concern for
42 the health of poultry ³. NH₃ produced through the decomposition of feces and urine by
43 microorganisms is a primary factor in the environmental pollution in poultry coops and can damage
44 the health of the respiratory system, eyes, paranasal sinuses, skin, and other organs ^{4,5}. The high
45 concentrations of NH₃ directly harm the immune function, health, and growth capability of poultry
46 ^{3,6,7}, gives rise to a variety of diseases and leads to economic losses. Currently, research on the
47 influence of NH₃ on poultry has focused on laying and broilers hens, while the effects of NH₃ on
48 geese have not been widely discussed ⁸. However, meat farming and laying geese in China are now
49 gradually transforming from the outdoor or semi-outdoor to the indoor type. Therefore, the hazards
50 associated with NH₃ are expected to pose a severe problem for meat geese production. Due to the
51 severe implications of NH₃ on poultry health, there is a need to provide a beneficial and stable
52 environment for poultry. Which is suitable for the complex nature of the NH₃ for its modelling and
53 prediction.

54

55 In recent years, several intelligent algorithms have been proposed for the prediction of NH₃ levels ⁹⁻¹³.
56 Artificial neural network (ANN), support vector machine (SVR), and decision tree (DT) are useful
57 tools and have been widely used for solving complex prediction problems. However, DT often faces
58 the over-fitting issue, so that it performs well on the training data set but not on the test set. SVM
59 uses the quadratic programming approach to measure the supporting vector making it difficult in
60 large-scale training sets to implement and make its output heavily dependent on the choice of
61 various hyperparameters. Instead, ANN has a strong generalization ability. Still, it can easily fall into
62 a partial solution ¹⁴. Ensemble learning integrates the prediction of several foundation estimators
63 established with a given learning algorithm to enhance the generalization performance over a single
64 estimator, and has become a hotspot in the prediction field and has been successfully applied in
65 some fields ¹⁵. Random forest (RF) is a representative ensemble learning method. Compared with the
66 methods mentioned above, RF with fewer hyper-parameters seldom over-fits and is relatively robust
67 to outliers and noise ¹⁶⁻¹⁸.

68

69 Many studies have indicated that changes in NH₃ are related to temperature, humidity and other
70 environmental factors.¹² used an adaptive neuron fuzzy inference system (ANFIS) to predict NH₃
71 using indoor relative humidity, indoor temperature, pig temperature and other indicators.¹⁹
72 predicted NH₃ based on genetic algorithm (GA) and optimized backpropagation neural network.
73 Temperature, relative humidity, carbon dioxide, total suspended particulates, solar radiation and
74 atmospheric pressure, have been usually used as prediction indicators²⁰. While these models select
75 some of the indicators as inputs to predict Nh₃, few studies have considered the correlation between
76 each feature and NH₃, and the methods used in these studies display several shortcomings. For
77 example, using the vast data directly in intelligent models without feature selection not only
78 increases the training time but also increases the risk of over-fitting.

79

80 Moreover, the contribution of the algorithm optimization and model combination may be lower than
81 that of screening for good prediction indicators¹⁸. Thus, feature selection is necessary. Extraction and
82 selection are common operations in feature engineering. Feature selection is better than feature
83 extraction with respect to readability and understandability and will not alter the primitive feature
84 data²¹. Recently, mutual information-based algorithms (MI) play an increasingly significant role in
85 data mining and machine learning. These methods have good non-linear and linear processing
86 capability for considering the relation of diverse sets of features^{22,23}. However, mutual information
87 algorithms ignore the influence of the proportions of labels on the correlation degree between
88 features and label sets²⁴. The Relief algorithm presented by²⁵ was initially used for two-category
89 problems. To adapt to multiple category problems,²⁶ proposed the relief-F algorithm based on Relief
90 to deal with multilabel data and regression problem, which is a widely adaptable filter-based feature
91 evaluating algorithm²¹. In this study, the fusion of mutual information and relief-F was proposed to
92 improve the capability of feature selection and bring a feature selection that is more accurate.

93

94 In this paper, we investigate the NH₃ prediction. Our target is to accurately forecast NH₃ levels using
95 the data from an intelligent geese house Internet of Things system. To overcome this challenging
96 problem, we design an RF underpinned framework that effectively predicts the NH₃ level. Although
97 RF is a promising approach, two challenges must be addressed. The first is that unrelated and
98 redundant features give rise to a high cost of the RF training process and decrease the prediction

99 accuracy. The second is the tuning of the parameters. There are three hyper-parameters: number of a
100 tree, size of sampling subsets and the minimum number of samples required to split an internal
101 node. These hyper-parameters affect the performance of RF, and currently, there is no consensus
102 regarding how their values should be set.

103

104 To address the challenges mentioned above, this paper combines the feature selector with the RF and
105 designs a new hybrid prediction model for predicting NH₃ in geese houses. First, the combined
106 feature selector (CFS) that combines the MI and relief-F evaluates the importance of prediction
107 indicators. Then, the parameter M is controlled to eliminate the features that have low feature
108 importance. Finally, the parameters of RF are used by CVGS to build a model for predicting NH₃ in
109 geese houses. The hybrid model makes full use of the CFS and is highly suitable for the selection of
110 the prediction indicators in this study. Comparison with other models horizontally and vertically
111 using empirical results shows that the prediction results of this model can provide technical support
112 for the precise management and control of intensive aquaculture environment of waterfowl.

113

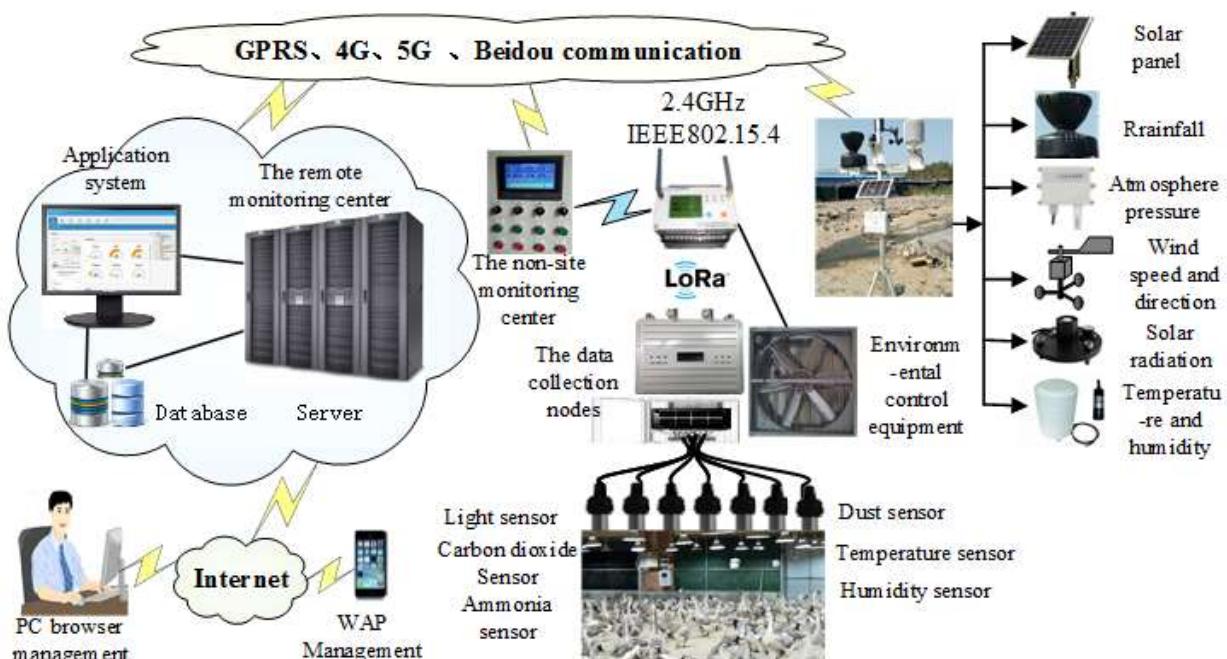
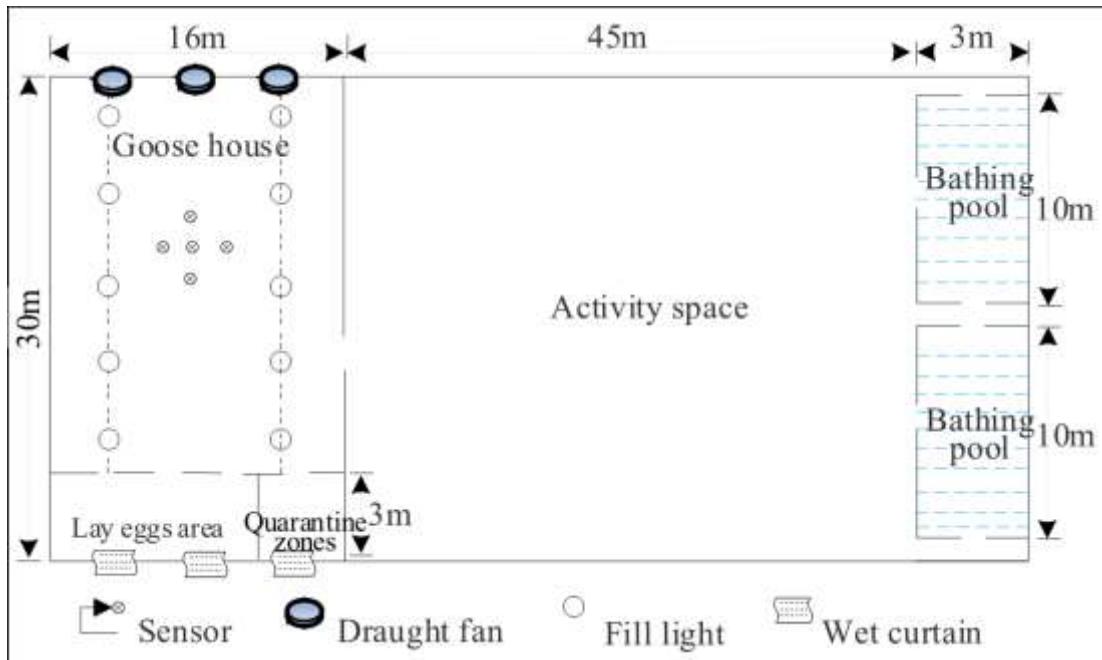
114 Materials and Methods

115 Study area and data source

116 The data used in this study were obtained from a waterfowl breeding farm in Haifeng County
117 (23 °05'N, 115 °19'E) in Shanwei city, China. With an area of approximately 53.3 hm², the farm is a
118 multifunctional integrated aquaculture base integrating waterfowl breeding, seeding breeding and
119 intensive aquaculture. In this experiment, we housed animals (stone goose) in the area including
120 poultry house (25 × 16 m²), playground (25 × 30 m²) and a swimming pool (20 × 3 × 1 m³). Fans,
121 control equipment (temperature and light), and various sensors were installed for online monitoring
122 of aquatic environment parameters in the geese houses(fig. 1).

123

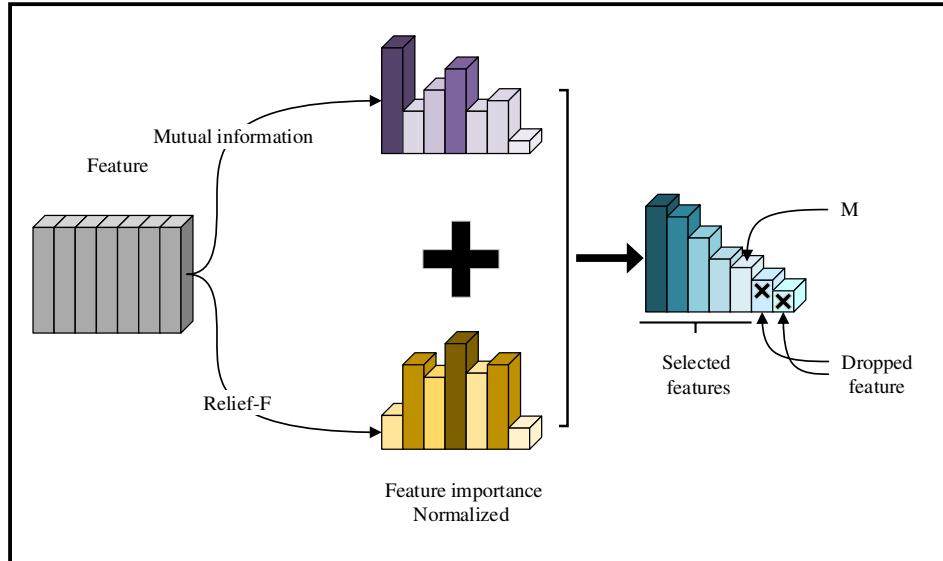
124 Because geese house environment parameters are mainly affected by physical and chemical factors.
125 We developed an IoT system(Fig2) to monitor temperature, humidity, carbon dioxide(CO₂), total
126 suspended particulates (TPS), NH₃ of the houses and temperature, humidity, atmospheric pressure,
127 and solar radiation of the surrounding environment.



135 Combined Feature Selector

136 To fully evaluate feature importance between each feature and target, MI and relief-F are both used
137 in CFS. These processes compute the feature importance of each feature to decide whether a feature
138 is eliminated or preserved. The parameter M is the threshold. We can eliminate the feature with low
139 feature importance by controlling M . In this section, we will introduce this module. Fig. 3 shows the

140 processes of CFS.



141
142 Fig.3. The structure of CFS.
143
144 **Relief-F**
145 The relief algorithm is an effective filtered feature selection method proposed by Kira and Rendell ²⁵.
146 The relief algorithm is initially limited to the classification of two types of data so the relief-F
147 algorithm, which is later extended by Kononeill can solve the multi-class and regression problems ²⁶.
148 This algorithm is a weighted algorithm that assigns weights to each feature according to the
149 relevance of the target. The larger the feature weight, the higher the contribution of the feature, and
150 vice versa, the lower the feature classification contribution.

151
152 Relief-F algorithm estimates feature weight according to the degree of distinguishing samples that
153 are close to each other based on the value of the feature. For that purpose, the relief-F algorithm
154 randomly selects a simple X_i (X_i has classified p) from training set D which has $|y|$ class. Then
155 searches for k of its nearest neighbours from the same class, called near-hit $X_{i,nh}$, and also k nearest
156 neighbours from each of the different classes, called near-miss $X_{i,j,nm}(j=1,2,\dots, |y|; j \neq p)$, then the weight
157 of feature L can be computed as follows:

$$158 \quad \delta^L = \sum_i -\text{diff}(X_i^L, X_{i,nh}^L)^2 + \sum_{j \neq p} (q_j \times \text{diff}(X_i^L, X_{i,j,nm}^L)^2) \quad (1)$$

159 Where q_j is the proportion of class j samples in data set D , $\text{diff}(a^j, b^j)$ denoted distance between simple
160 a and b in feature j , $\text{diff}()$ defined as:

161
$$diff(a^j, b^j) = \begin{cases} \frac{|a^j - b^j|}{\max(j) - \min(j)}, & \text{if } j \text{ is continuous} \\ 0, & \text{if } j \text{ is discrete and } a^j \neq b^j \\ 1, & \text{if } j \text{ is discrete and } a^j = b^j \end{cases} \quad (2)$$

162

163 **Mutual Information Estimator**

164 Among the estimates of independence between random variables, MI is selected based on its
 165 information theory background. $MI(X, Y)$ between the two random variables X and Y is defined by
 166 the common information found in two variables with a joint probability distribution $P(X, Y)$. $MI(X, Y)$
 167 computes the degree of correlation between vector X and target vector Y and is given by:

168
$$MI(X, Y) = \int_Y \int_X P(x, y) \log \left(\frac{P(x, y)}{P(x)P(y)} \right) dx dy \quad (3)$$

169

170 where $P(x, y)$ is the probability density function of random variable $Z=(X, Y)$, $P(x)$ and $P(y)$ represents
 171 the marginal probability density function of X and Y , respectively. In fact, since we usually do not
 172 know $P(x, y)$ in advance, some methods should be used to estimate $MI(X, Y)$. K-nearest neighbour is a
 173 non-parametric method that has been confirmed as a useful method in MI estimation²⁷. For K-
 174 nearest neighbour, we use Euclidean distance as a distance metric, and the maximum norm for the
 175 space $Z=(X, Y)$ written as:

176
$$\|z - z'\| = \max \{ \|x - x'\|, \|y - y'\| \} \quad (4)$$

177

178 Let $\kappa(i)/2$ express the distance from z_i to its k^{th} neighbour, and $\kappa_x(i)/2$ and $\kappa_y(i)$ express the
 179 distance between the same points projected into the X and Y subspaces. It is clear that:

180
$$\kappa(i) = \max \{ \kappa_x(i), \kappa_y(i) \} \quad (5)$$

181

182 Then, we denote by $\eta_x(i)$ the number of points for which the distance from x_i is strictly less than $\kappa(i)$
 183 and by $\eta_y(i)$ is the number of points for which the distance from y_i is strictly less than $\kappa(i)$. We note
 184 that $\kappa(i)$ not a fixed value, and $\eta_x(i)$ and $\eta_y(i)$ is also not fixed. We denote by $\langle \cdot \cdot \cdot \rangle$ both all $i \in (1, \dots, N)$
 185 and all realizations of the random samples:

186

$$\langle \dots \rangle = \frac{\sum_{i=1}^N E[\dots(i)]}{N} \quad (6)$$

187

188 The estimate for MI by K-nearest neighbour is then:

189

$$I(X, Y) = \psi(k) - \langle \psi(\eta_x + 1) + \psi(\eta_y + 1) \rangle + \psi(N) \quad (7)$$

190

191 In equation [], $\psi(\cdot)$ is the di-gamma function and satisfies the recursion $\psi(x+1)=\psi(x)+1/x$ and $\psi(1)=-C$,
 192 where $C=0.5772156..$ is the Euler-Mascheroni constant. If MI is equal to zero, the two random
 193 variables are independent and higher MI values mean higher dependency.

194 **Random Forest.**

195 The RF algorithm is a non-linear ensemble model that establishes and averages a large number of
 196 random distribution DT for regression or classification tasks ²⁸. A DT or classification and regression
 197 tree (CART) that construct the RF is a non-parametric model. According to the complexity of the
 198 input data, the tree grows in the process of the learning. Decision nodes and leaf nodes are the main
 199 components of DT. Each input sample is estimated by a test function of decision nodes and passed to
 200 different branches according to the features of the sample. Let we denote by $X=\{x_1, x_2, x_3, \dots, x_n\}$ the
 201 input vector with n features, Y the output scalar and D_m the training set with m observations which
 202 can be written as follows:

203

$$D_m = \{(X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \dots, (X_m, Y_m)\}, X \in R^n, Y \in R \quad (8)$$

204

205 At each node, the input data are split by a specific algorithm in the process of training to optimize the
 206 parameters of the split function to the fit data set D_n . In the first step, the DT must be optimally split
 207 among all of the variables. The splitting procedure begins at the root node, and each node uses its
 208 own split function for the new input X . This operation is recursive until a leaf node appears. The tree
 209 stops growing either when the maximum number of levels is reached or when the observation
 210 number of a node is less than a predefined number. The result of the DT learning process is a
 211 prediction function $\hat{T}(D_m, X)$ generated over D_m .

212

213 RF regression model can offer powerful prediction ability and is an extension of the DT. The main
 214 characteristics of RF include bootstrap resampling and random feature subset. An RF is an ensemble
 215 of P DT $\hat{T}(D_m^1, X), \hat{T}(D_m^2, X), \dots, \hat{T}(D_m^P, X)$. Here, $(D_m^1, D_m^2, \dots, D_m^P)$ are the bootstrap sample obtained
 216 by random sampling of m observations with replacement from D_m , where each observation has
 217 $1/m$ probability of being drawn. This sample process is known as bootstrap resampling. During the
 218 splitting of each node, only a small part of n feature are randomly selected instead of all features;
 219 this is known as random feature selection. The ensemble learning result P output $\hat{Y}_1 = \hat{T}(D_m^1, X)$,
 220 $\hat{Y}_2 = \hat{T}(D_m^2, X), \dots, \hat{Y}_P = \hat{T}(D_m^P, X)$. Then, the final estimation output \hat{Y} is the average of P output, that
 221 is described as follows:

$$\hat{Y} = \frac{1}{P} \sum_{i=1}^P \hat{Y}_i = \frac{1}{P} \sum_{i=1}^P \hat{T}(D_m^i, X) \quad (9)$$

223
 224 where \hat{Y}_i is the output of i^{th} DT, $i = 1, 2, \dots, P$. The framework of RF regression is illustrated in
 225 Fig.4, and its training process can be summarized as:

- 226 1. Obtain bootstrap samples from the training data set by bootstrap resampling
 227 2. Generate a regression DT by full use of the bootstrap sample drawn in step 1 with the following
 228 modification: at each node, select the optimal split among a random subset sampled in input
 229 variables($mtry$) instead of all of them.
 230 3. Repeat steps 1 and 2 until P DT tree is generated.
 231 4. Aggregating the output of P trees by an average method to forecast unknown data.

232 *2.4 Cross-validation grid search*
 233

234 According to the previous section, we notice that the RF algorithm has two important parameters:

- 235 • P : the number of decision trees that are the base estimators of RF.
 236 • $mtry$: the size of the random feature subset.

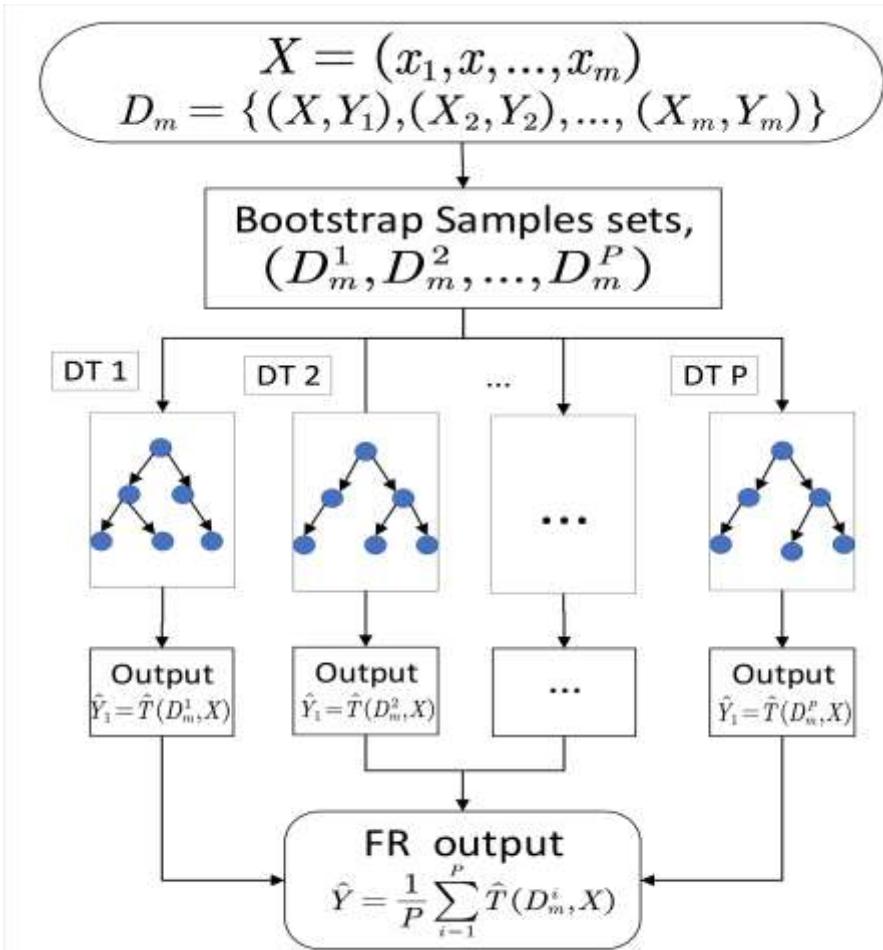
237 Generally, a variance of RF decreases as P grows. More accurate predictions are likely to be
 238 obtained by choosing a large number of trees, but there is no common setting for P ²⁹. Additionally,
 239 $mtry$ also a sensitive parameter, and increasing $mtry$ can improves the intensity of each DT but the

240 relation among DT will also be increased. This means that the total strength of RF may be decreasing.
241 Therefore, it is necessary to optimize the parameters of RF and select the optimal RF parameters.

242
243 The grid search algorithm that is currently the most widely used method for parameter optimization
244 is a highly suitable model with fewer hyper-parameters. GS exhaustively generates candidates from
245 a grid of parameter values specified by the user parameters, and then trains each candidate set of
246 parameters and marks the score of the model, finally obtaining the optimal combination of
247 parameters. GS optimizes all of the parameters of the model to guarantee that the given best
248 parameter combination is the optimal global solution in the pre-setting grid.

249
250 At the same time, learning the parameters of a model and testing it on the same data set is a mistake,
251 because in this case, the model that merely repeats the samples that it has just learned will have a
252 high score. To avoid this, we combine K-fold cross-validation (CV) and grid search (GS). The score is
253 evaluated based on the mean square error (MSE) average value given from K iterations. Finally, the
254 optimal parameter combination with the lowest MSE is obtained. The specific steps are described as
255 follows:

256



257

258 **Fig.4. The framework of random forest regression.**

259

- 260 • Step 1: Partitioning the train data into equal-size k sets.
 261 • Step 2: Setting the scale of all of the parameters and exhaustively generating candidates from the
 262 parameter space.
 263 • Step 3: A model with parameters combination is trained using $K - 1$ of the folds as the training
 264 set. The MSE of the model is computed on the remaining part of the data.
 265 • Step 4: Each of the K folds followed the step 3 procedure.
 266 • Step 5: The score measured by CVGS is the average of the values computed in the step 4. Then,
 267 the parameter combination with the lowest MSE is identified as the optimal.

268 **Hybrid prediction model based on CFS-CVGS-RF.**

269 This paper proposes a CFS-CVGS-RF model to predict NH₃ in a geese house. The methodology for
 270 conducting this model is shown in Fig.5. The implementation process for NH₃ prediction based on
 271 CFS-CVGS-RF can be described as follows:

272

- 273 • Step 1: Data normalized processing. Different data of geese house environment has other units
 274 and dimensions, using the original data directly will make the model complex and will also
 275 decrease the performance of prediction. To address this problem, we use the normalization that
 276 can eliminate the difference between the data units and dimensions and facilitates the study of
 277 the correlation between environmental factors. The normalized process method is described by:

278
$$y'_i = \frac{y_i - \bar{y}}{y_{\max} - y_{\min}} \quad (10)$$

279 where y'_i are the normalized data, y_{\max} and y_{\min} are the max and min values of the original
 280 data, and y_i and \bar{y} are the original data and its mean.

281

- 282 • Step 2: Feature selection based on CFS. CFS selects the features that are strongly related to NH_3
 283 and eliminates the low-importance factors. The remaining features are used as the input to the
 284 regression model. CFS reduces the dimension of input and solves the problem of information
 285 redundancy. In the CFS process, we first compute linear and non-linear correlation strengths
 286 between each environmental factor and target(NH_3) by using relief-F and MI, respectively. The
 287 final feature importance is the sum of two different dimensions features important values after
 288 normalization. Then, threshold M is set to the screen factor with feature importance lower than
 289 M .

290 •

- 291 • Step 3: CVGS-RF modelling. RF has two key parameters, namely, the number of DT P and size
 292 of the random feature subset $mtry$. To find the optimal values of these two parameters, we
 293 adopt the CVGS method. In the CVGS process, we first establish the grid coordinates of the
 294 parameters. In this paper, combining the RF related literature ²⁹ and the experimental
 295 parameters wave motion, we set $P = [2,1000]$ and $mtry = [2,4]$. Then, the data set is divided
 296 into K subsets, where 10-fold ($k = 10$) cross-validation is considered to be better ³⁰. After k
 297 parallel operations, each parameters' combinations have k MSE. According to the mean of each
 298 calculation results, we select the parameter combination with minimum average MSE as the
 299 optimal parameters and establish the RF model using these values.

- 300 • Step 4: Result output denormalization. Denormalize the output to obtain the results in the
 301 normal dimension. The denormalize process is described by:

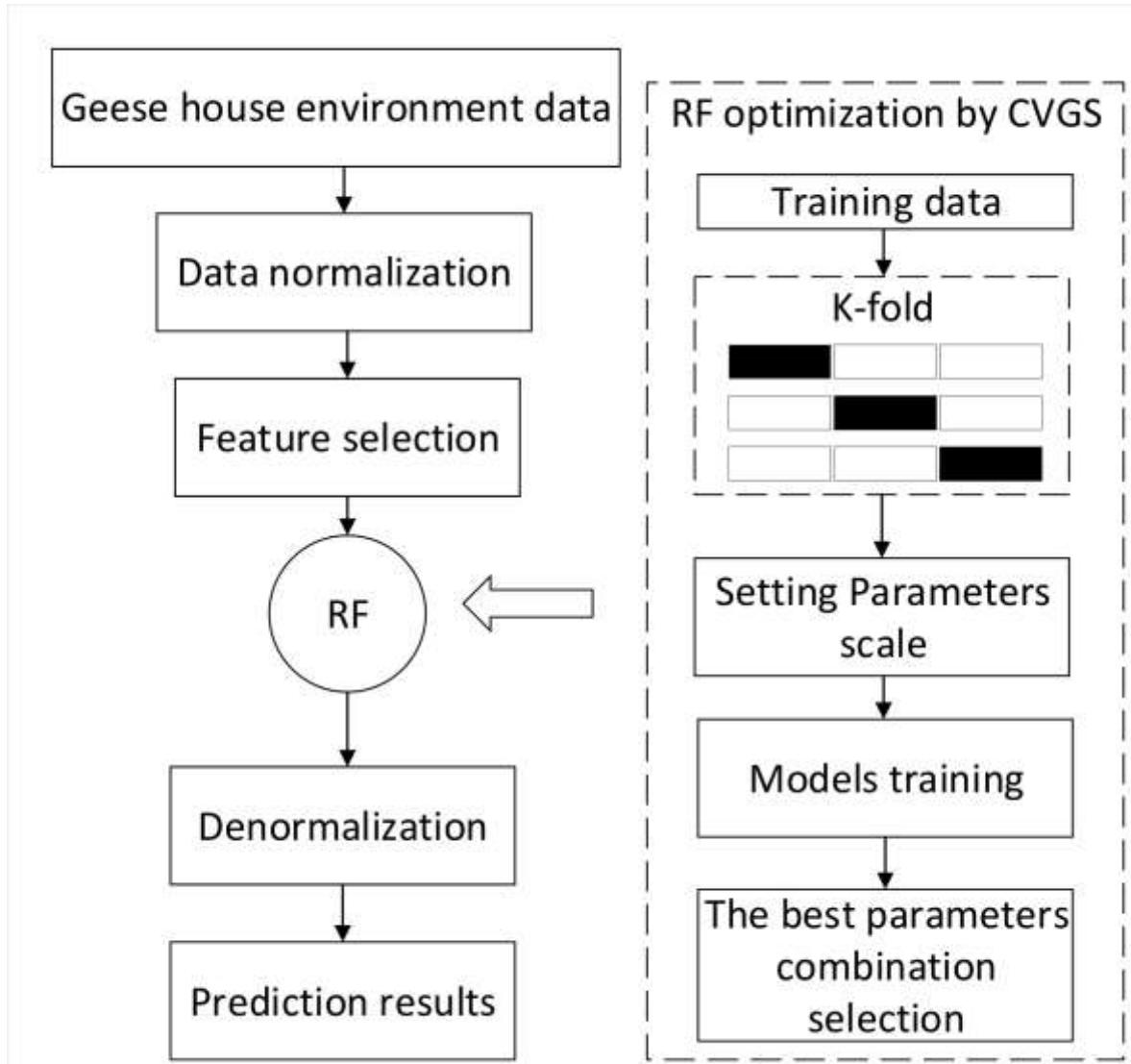
302
$$y_i = y'_i * (y_{\max} - y_{\min}) + \bar{y} \quad (11)$$

303

304 where y_i' is the denormalized data, y_{max} and y_{min} are max and min value of original data, and y_i'
305 and \bar{y} are the original data and its mean.

306

307 In this article, a simulation model was built in order to validate the performance of the proposed
308 method. The NH₃ was tested in a high-density geese culture farm from September 10st to September
309 27st, 2019 in Sanwei city, Guangdong province. Fig.6 shown the feature importance computed by CFS
310 for each factor, and the performances of different threshold M are given in Table 1. The P and $mtry$
311 combinational parameters of optimal CVGS-RF model are 500 and 3, respectively, and the MSE
312 values for different parameter combinations are shown in Fig.7. It is observed that after growing 100
313 trees, the MSE value decreases very slowly and apparently converges after 500 trees, with $mtry=3$
314 resulting in the lowest MSE. Fig.8 shows the NH₃ series prediction result obtained by the combined
315 model based on CFS-CVGS-RF. The performance estimation statistics of the testing are given in Table
316 2. Fig.9 shows the residual distribution of the proposed model. The results of NH₃ prediction
317 demonstrate the robustness and effectiveness of CFS-CVGS-RF.



318

319

Fig.5. The schematic of the proposed methodology.

320

321 **Results and discussions**

322 **Experimental environment and settings**

323 The geese house environment data considered in this investigation include the data obtained at the
 324 intervals of twenty minutes from September 10 to September 27, 2019. Seventy-two sets of data
 325 collected per day yield a total of 1296 observations samples. For a model generation, the first 864 sets
 326 of the data were used for model training, and the remaining 432 sets were used as the testing data to
 327 estimate the prediction performance of the constructed model.

328

329 To confirm the superiority of the framework constructed in this research, we designed a script with
 330 Python 2.7vision following the system framework described in Section 3. The script is executed on

331 the Win10 operation system with Intel Core i5, 4 GB RAM, and 500G hard disk. We use the optimal
332 parameter combination ($P=500$, $mtry=3$) for the prediction model of NH₃ by the CVGS algorithm.
333 According to Fig.6 and Table2, the used features include outdoor humidity, TPS, CO₂ and indoor
334 temperature as the optimal input for the model to forecast NH₃ in this paper.

335

336 Performance criteria

337 For a reasonable evaluation of each prediction model, three commonly used error standards are
338 proposed to measure the prediction accuracy of the model, including MSE, RMSE and MAPE. The
339 relevant calculation formulae are:

340

$$341 \quad MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (12)$$

$$342 \quad RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} \quad (13)$$

$$343 \quad MaxAPE = Max \left| \frac{y_i - \hat{y}_i}{\hat{y}_i} \right| \times 100\% \quad (14)$$

344

Results and analysis

345 Following the method introduced in this article, we use relief-F to compute the linear relation
346 strength between each factor and NH₃, and we use MI rather than relief-F to compute the non-linear
347 relation strength. Fig.6 shows the normalized results of relief-F and MI, and it is observed that
348 outdoor humidity has the highest relation and solar radiation has the lowest relation with NH₃. To
349 find the optimal input feature combination, the threshold M was changed from 0 to 2. The sensitivity
350 test results of M were recorded in Table 1, and it is observed that the combination of outdoor
351 humidity, TPS, CO₂ and indoor temperature features was the optimal input combination. Fig.7
352 illuminates the change of the fitness value, with the three convergence curves showing the best fit for
353 RF. Three fitness curves show that after growing 100 trees, the MSE decreases very slowly and
354 apparently converges after 500 trees, and $mtry=3$ results in the lowest MSE.

355

356 To analyze and compare prediction performance, two types of comparison were designed: (1)
357 horizontal comparison between the model with CFS and the model without CFS and (2) vertical
358 comparison between the models used in this paper with other parallel models. The horizontal
359 comparison includes the CVGS-RF model, and vertical comparison consists of the benchmarks used

360 in this DT, support vector machine and back propagation neural network (BPNN). To verify the
361 performance forecasted by the models in this paper, these models used data sets and predict the NH₃
362 content of the last 72 test sets corresponding to the last 24 hours. Fig.8 and fig.9 show the prediction
363 curves and the error bar plot.

364

365 It is observed from the figures that the prediction curve of the CFS-CVGS-RF model is closer to the
366 original value than the prediction curves of the other four models, and has better accuracy than the
367 other four models. It can be seen from the prediction error box plot (fig.10) that the CFS-CVGS-RF
368 model has smaller error fluctuations than the other model, and the CVGS-RF model using feature
369 combination after feature selection by CFS has smaller forecast error fluctuations than CVGS-RF
370 without feature selection, indicating that CFS is effective.

371

372 **Table.1. The sumulation results of CFS.**

M	Accuracy	Running time(second)	Eliminated feature
0	82.0%	6.817	
...
1.0	88.1%	6.489	Solar radiation
1.2	94.3%	6.286	Atmospheric pressure
1.4	95.6%	5.850	Outdoor temperature
1.6	78.5%	5.194	Indoor humidity
1.8	16.1%	4.945	TPS
2.0	NaN	NaN	Indoor temperature
			CO ₂
			Outdoor humidity

373

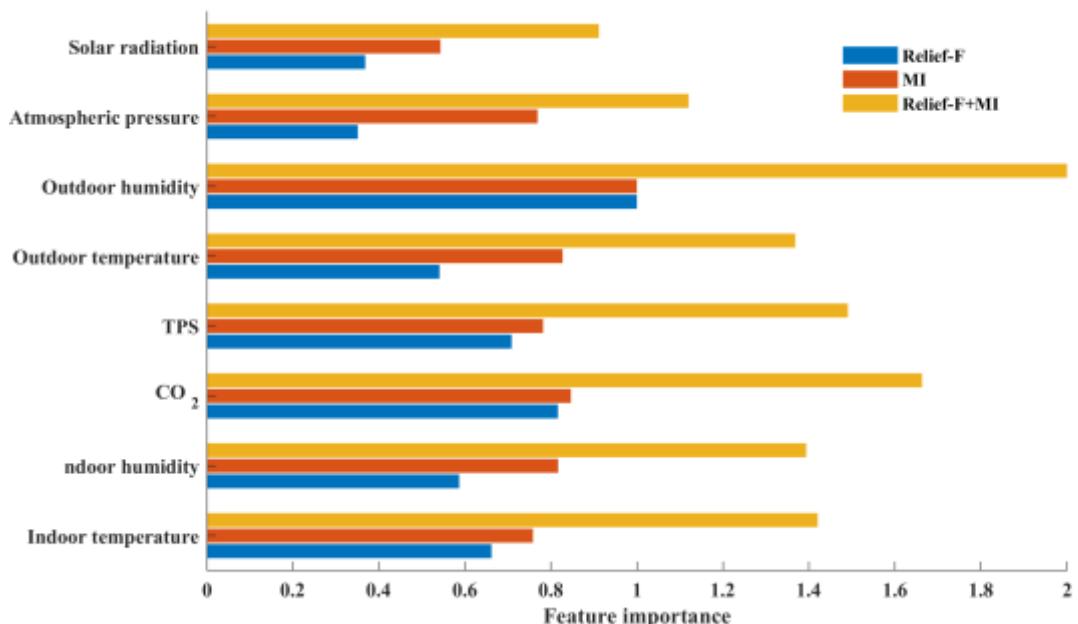
374 For a more accurate comparison of the performance of five models, this article computes MSE, RMSE
375 and MaxMAPE of those models for which the details are shown in Table 2. The MSE, RMSE and
376 MAPE of CFS-GS-RF and GS-RF were 0.5072, 0.6583, 2.88% and 1.2658, 1.2851, 9.10%, respectively.
377 These values are the best estimation indexes among the five models. Fig.8 describes the prediction
378 residual distribution condition of five models, and it is observed that CFS-GS-RF has fewer outliers
379 and overall error closer to zero for the residual error compared with other models. This mean that
380 CFS-GS-RF has more stability for forecasting results and is more suitable for the prediction of NH₃.

381

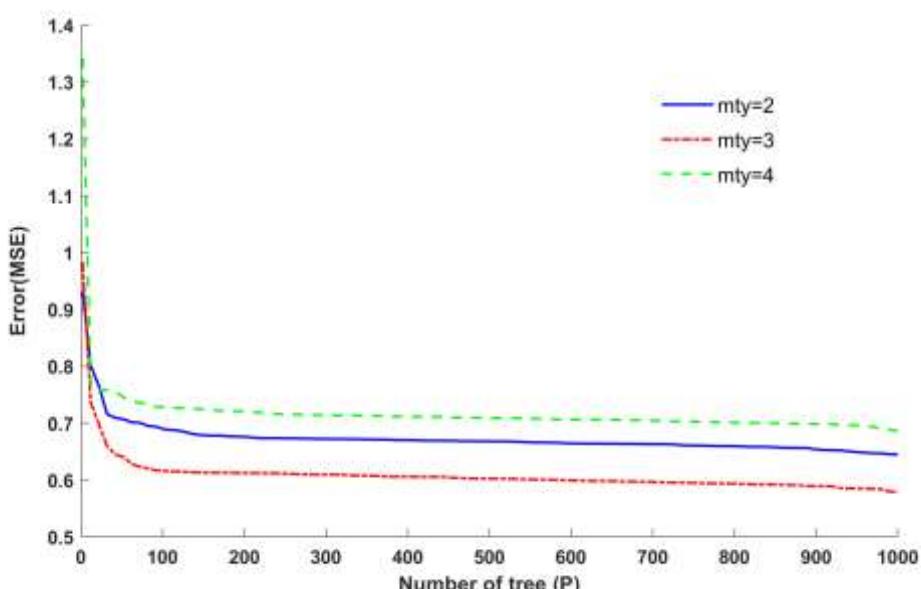
382 **Conclusions**

383 This paper proposes a novel NH₃ prediction hybrid model (CFS-CVGS-RF). The CFS-CVGS-RF

384 model consists of a combination of four methods: relief-F, mutual information, K-fold cross-
 385 validation grid search optimization algorithm and random forest. Using actual experimental geese
 386 house environment data from a monitored aquaculture factory farm in Sanwei city, China. Results
 387 clearly show that the proposed hybrid method CFS-CVGS-RF has better forecasting performance
 388 than CVGS-RF, DT, SVM and BPNN, as measured by MSE, RMSE and MaxMAPE. Furthermore,
 389 CFS-CVGS-RF can effectively consider the linear and non-linear relations between the input features
 390 and the target, reduce redundant information and improve the prediction performance of the model
 391 by screening the unrelated or low-relation features.



392
 393 Fig.6. Feature importance computed by Relief-F and MI.



394
 395 Fig.7. The results of K-fold cross-validation grid search.

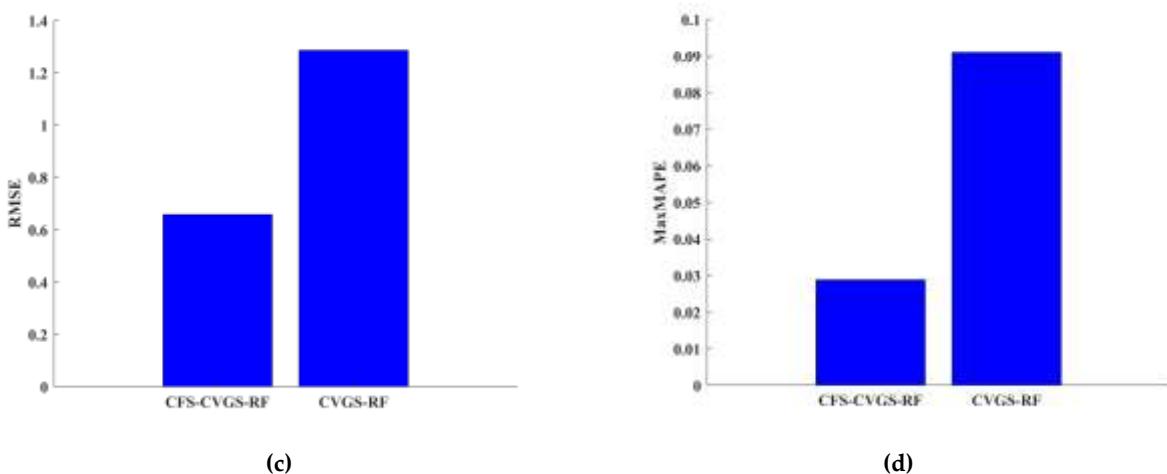
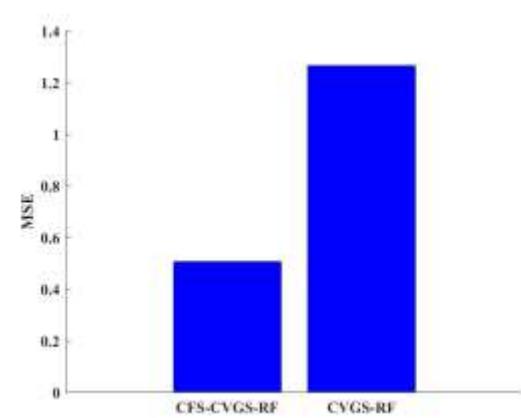
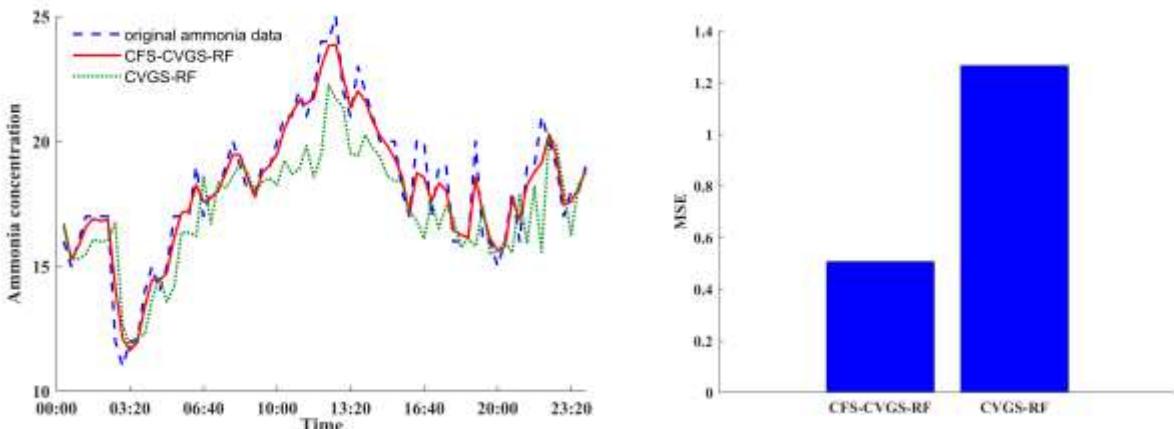
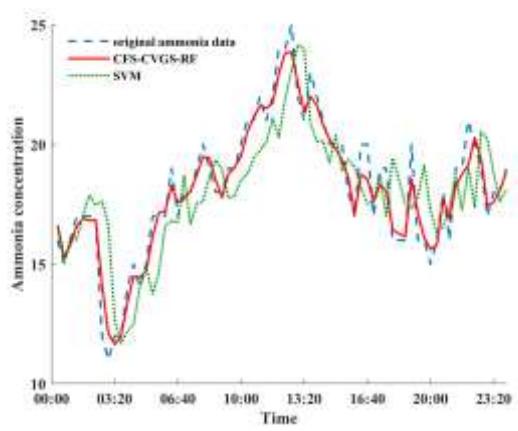
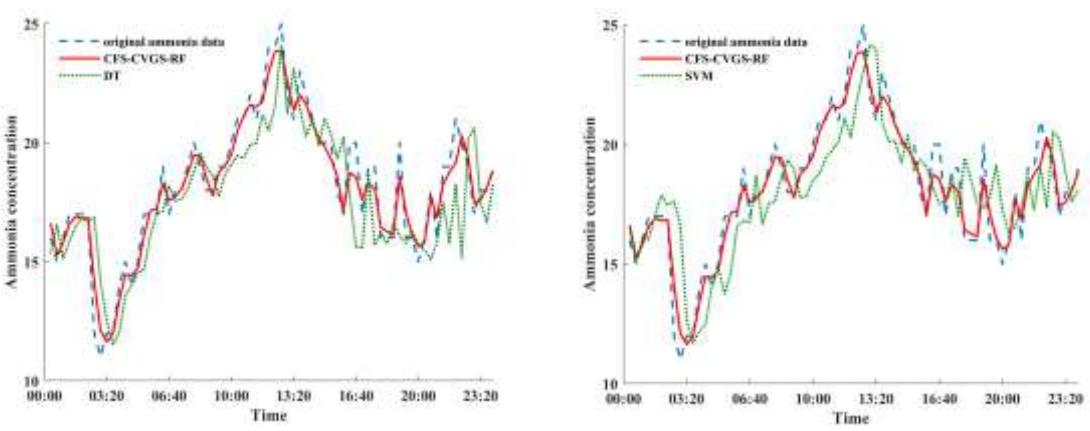


Fig.8. Horizontal comparison results.



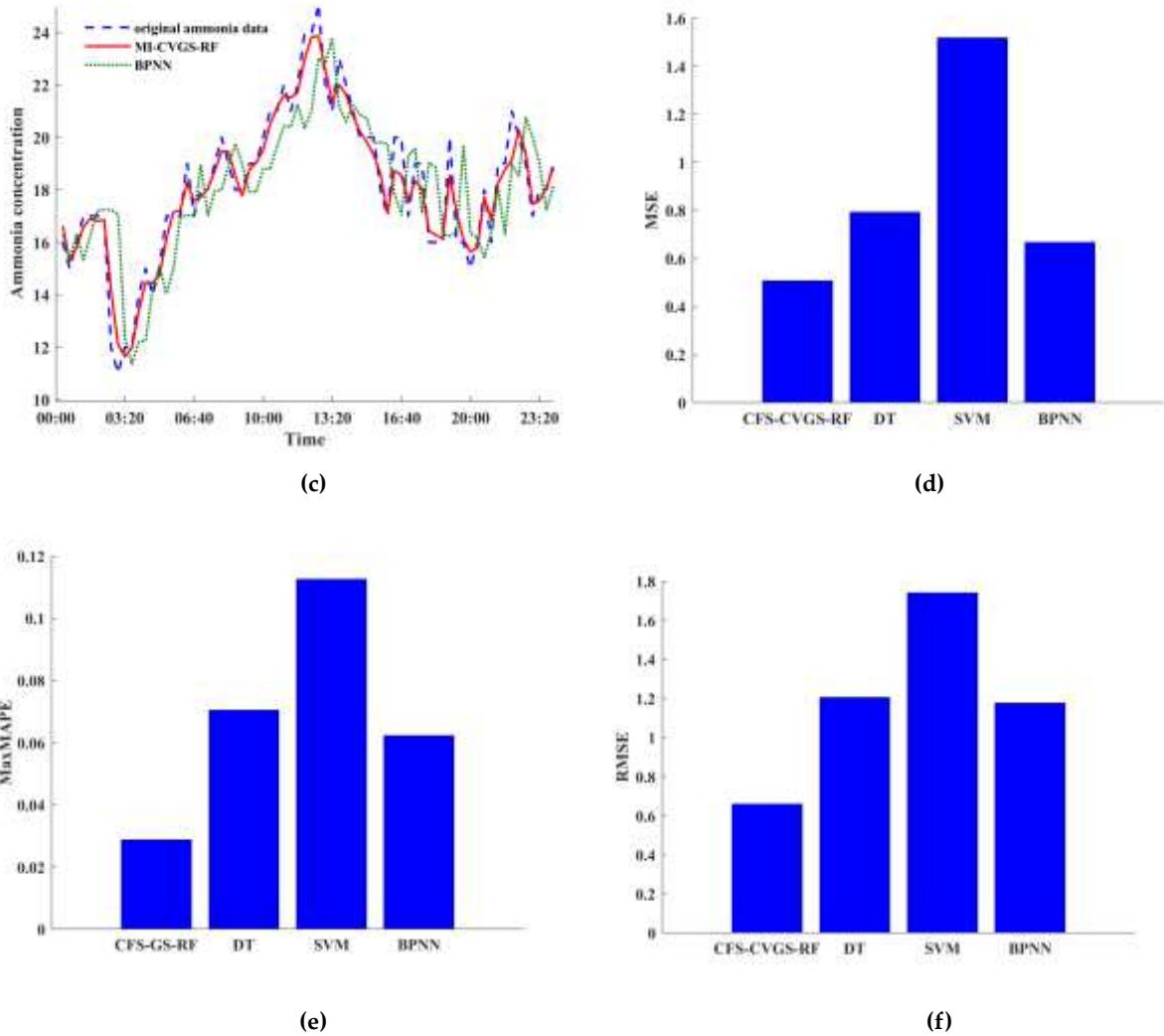
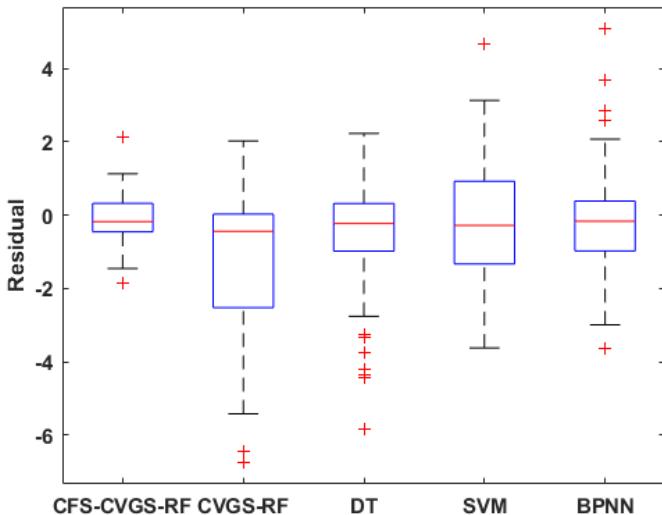


Fig.9. The vertical comparison results.

Table.2. Comparison of NH₃ prediction results.

Model	CFS-CVGS-RF	CVGS-RF	DT	SVM	BPNN
MSE	0.5072	1.2658	0.7922	1.5179	0.6667
RMSE	0.6583	1.2851	1.2047	1.7400	1.1764
MaxMAPE(%)	2.88	9.10	7.04	11.25	6.23

Our study has several limitations that require further research. First, predicting NH₃ is a very complex issue that is influenced by many factors; however, due to equipment limitations, we cannot monitor more related factors that may strongly influence NH₃. Second, concerning experimental time, in the future, we plan to collect data in other months to verify whether the proposed model is useful for a different season. Finally, we plan to investigate the use of other ensemble strategies instead of a simple averaging method to improve the RF, such as weighted averaging, weighted voting.



419

420 Fig.10. Boxplot of Residual error in different models.

421

422 **Data availability**

423 Data or code presented in this study are available on request from the corresponding author.

424

425 **Reference**

- 426 1. CSB, C. S. B. Statistical Communique of the People's Republic of China on the 2019
427 National Economic and Social Development. (2020).
- 428 2. Zhao, Q., Boomer, G. S. & Kendall, W. L. The non-linear, interactive effects of
429 population density and climate drive the geographical patterns of waterfowl survival.
430 *Biol. Conserv.* **221**, 1–9 (2018).
- 431 3. Wei, F. X. *et al.* Ammonia concentration and relative humidity in poultry houses affect
432 the immune response of broilers. *Genet Mol Res* **14**, 3160–3169 (2015).
- 433 4. Kearney, G. D., Shaw, R., Prentice, M. & Tutor-Marcom, R. Evaluation of respiratory
434 symptoms and respiratory protection behavior among poultry workers in small
435 farming operations. *J. Agromedicine* **19**, 162–170 (2014).
- 436 5. Nemer, M. *et al.* Airway inflammation and ammonia exposure among female

- 437 Palestinian hairdressers: a cross-sectional study. *Occup. Environ. Med.* **72**, 428–434
438 (2015).
- 439 6. Xiong, Y., Tang, X., Meng, Q. & Zhang, H. Differential expression analysis of the
440 broiler tracheal proteins responsible for the immune response and muscle contraction
441 induced by high concentration of ammonia using iTRAQ-coupled 2D LC-MS/MS. *Sci.
442 China Life Sci.* **59**, 1166–1176 (2016).
- 443 7. Soliman, E. S., Moawed, S. A. & Hassan, R. A. Influence of microclimatic ammonia
444 levels on productive performance of different broilers' breeds estimated with
445 univariate and multivariate approaches. *Vet. world* **10**, 880 (2017).
- 446 8. Tao, Z. *et al.* Effects of ammonia on intestinal microflora and productive performance
447 of laying ducks. *Poult. Sci.* **98**, 1947–1959 (2019).
- 448 9. Bai, Y., Li, Y., Wang, X., Xie, J. & Li, C. Air pollutants concentrations forecasting using
449 back propagation neural network based on wavelet decomposition with
450 meteorological conditions. *Atmos. Pollut. Res.* **7**, 557–566 (2016).
- 451 10. Lim, Y., Moon, Y.-S. & Kim, T.-W. Artificial neural network approach for prediction of
452 ammonia emission from field-applied manure and relative significance assessment of
453 ammonia emission factors. *Eur. J. Agron.* **26**, 425–434 (2007).
- 454 11. Qiao, J., Quan, L. & Yang, C. Design of modeling error PDF based fuzzy neural
455 network for effluent ammonia nitrogen prediction. *Appl. Soft Comput.* 106239 (2020).
- 456 12. Xie, Q., Ni, J. & Su, Z. A prediction model of ammonia emission from a fattening pig
457 room based on the indoor concentration using adaptive neuro fuzzy inference system.
458 *J. Hazard. Mater.* **325**, 301–309 (2017).
- 459 13. Stamenkovic, L. J., Antanasijevic, D. Z., Ristic, M. D. J., Peric-Grujic, A. A. & Pocajt, V.

- 460 V. Modeling of methane emissions using artificial neural network approach. *J. Serbian
461 Chem. Soc.* **80**, 421–433 (2015).
- 462 14. Yu, H., Chen, Y., Hassan, S. G. & Li, D. Prediction of the temperature in a Chinese
463 solar greenhouse based on LSSVM optimized by improved PSO. *Comput. Electron.
464 Agric.* **122**, 94–102 (2016).
- 465 15. Barzegar, R., Fijani, E., Moghaddam, A. A. & Tziritis, E. Forecasting of groundwater
466 level fluctuations using ensemble hybrid multi-wavelet neural network-based models.
467 *Sci. Total Environ.* **599**, 20–31 (2017).
- 468 16. Qiu, X., Zhang, L., Nagaratnam Suganthan, P. & Amaratunga, G. A. J. Oblique random
469 forest ensemble via Least Square Estimation for time series forecasting. *Inf. Sci. (Ny)*.
470 (2017) doi:10.1016/j.ins.2017.08.060.
- 471 17. Rubal & Kumar, D. Evolving Differential evolution method with random forest for
472 prediction of Air Pollution. in *Procedia Computer Science* (2018).
473 doi:10.1016/j.procs.2018.05.094.
- 474 18. Wen, L. & Yuan, X. Forecasting CO₂ emissions in Chinas commercial department,
475 through BP neural network based on random forest and PSO. *Sci. Total Environ.* **718**,
476 137194 (2020).
- 477 19. Zhu, K., Wu, S. & Li, Q. Prediction Model for Piggery Ammonia Concentration Based
478 on Genetic Algorithm and Optimized BP Neural Network. *Metall. Min. Ind.* (2015).
- 479 20. Rong, L., Nielsen, P. V, Bjerg, B. & Zhang, G. Summary of best guidelines and
480 validation of CFD modeling in livestock buildings to ensure prediction quality.
481 *Comput. Electron. Agric.* **121**, 180–190 (2016).
- 482 21. Wang, K., Xu, C., Zhang, Y., Guo, S. & Zomaya, A. Y. Robust big data analytics for

- 483 electricity price forecasting in the smart grid. *IEEE Trans. Big Data* **5**, 34–45 (2017).
- 484 22. Sun, L., Wang, L., Ding, W., Qian, Y. & Xu, J. Feature Selection Using Fuzzy
485 Neighborhood Entropy-Based Uncertainty Measures for Fuzzy Neighborhood
486 Multigranulation Rough Sets. *IEEE Trans. Fuzzy Syst.* (2020)
487 doi:10.1109/tfuzz.2020.2989098.
- 488 23. Gonzalez-Lopez, J., Ventura, S. & Cano, A. Distributed multi-label feature selection
489 using individual mutual information measures. *Knowledge-Based Syst.* (2020)
490 doi:10.1016/j.knosys.2019.105052.
- 491 24. Sun, L., Yin, T., Ding, W., Qian, Y. & Xu, J. Multilabel feature selection using ML-
492 ReliefF and neighborhood mutual information for multilabel neighborhood decision
493 systems. *Inf. Sci. (Ny)*. (2020) doi:10.1016/j.ins.2020.05.102.
- 494 25. Kira, K. & Rendell, L. A. Feature selection problem: traditional methods and a new
495 algorithm. in *Proceedings Tenth National Conference on Artificial Intelligence* (1992).
- 496 26. Kononenko, I. Estimating attributes: Analysis and extensions of RELIEF. in *Lecture
497 Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and
498 Lecture Notes in Bioinformatics)* (1994). doi:10.1007/3-540-57868-4_57.
- 499 27. Kraskov, A., Stögbauer, H. & Grassberger, P. Estimating mutual information. *Phys.
500 Rev. E* **69**, 66138 (2004).
- 501 28. Breiman, L. Random forests. *Mach. Learn.* **45**, 5–32 (2001).
- 502 29. Probst, P. & Boulesteix, A.-L. To tune or not to tune the number of trees in random
503 forest. *J. Mach. Learn. Res.* **18**, 6673–6690 (2017).
- 504 30. Witten & Frank, I. H. Data mining. *Pract. Mach. Learn. Tools Tech. with Java
505 Implementations* **13**, 1 (2005).

506

507 **Funding**

508 This work was supported in part by the National Natural Science Foundation of China under Grant
509 61871475,61471-131,61571444, in part by the special project of laboratory construction of Guangzhou
510 Innovation Platform Construction Plan under Grant 201905010006, Guangzhou Innovation Plantform
511 Construction Plan under Grant 2017B0101260016, foundation for High-level Talents in Higher
512 Education of Guangdong Province under Grant 2017GCZX00014, 2016K-
513 ZDXM0013,2017KTSCX094,2018LM2168, and Beijing Natural Science Foundation under Grant
514 4182023.

515

516 **Contributions**

517 All authors contributed extensively to this manuscript. J.H. and S.H. contributed to
518 conceptualization, methodology, writing original draft preparation, software, visualization and
519 supervision. J.G. performed the formal analysis, H.W. and X.Z. contributed to the investigation, S.L.
520 provided the farm and data resources for this research. All authors have read and agreed to the
521 published version of the manuscript.

522

523 **Competing interests**

524 The authors declare no competing interests.