

Predicting tobacco pyrolysis based on chemical constituents and heating conditions using machine learning approaches

Hao Wei^a, Jiangkuan Xing^{a,*}, Kun Luo^a, Yuhan Peng^b,
Jianren Fan^a, Ke Zhang^c, Hui Wang^{b,*}

^a*State Key Laboratory of Clean Energy Utilization, Zhejiang University, Hangzhou,
310027, China*

^b*China Tobacco Zhejiang Industrial Co., Ltd., Hangzhou 310088, P. R. China*

^c*Zhengzhou Tobacco Research Institute of CNTC, Zhengzhou 450001, P. R. China*

Abstract

Tobacco is a special type of biomass that consists of complex chemical constituents. Currently, only global kinetic models have been developed for tobacco pyrolysis, but accurate kinetics considering the effects of the complex chemical constituents and heating conditions have not been well established. To this end, a general tobacco pyrolysis model was developed based on the complex chemical constituents and heating conditions using machine learning approaches. Specifically, chemical analysis and thermogravimetric analysis (TGA) of 49 tobacco samples under a wide range of heating rates were first conducted by experiments and then used to construct a database for the model development. Subsequently, the constructed database was divided into seen and unseen data-sets for the model development and evaluation. General pyrolysis models for single/multiple heating rates were developed from the seen data-set using an advanced machine learning approach, the Extremely Randomized Trees (Extra-Trees, ET). The performances of mod-

els were further evaluated on the unseen data-set through comparisons with the experimental data. The results showed that after feature selection based on Pearson correlation coefficient and hyper-parameters optimization, the trained models could accurately reproduce the tobacco pyrolysis behaviour on the unseen data with $R^2 > 0.967$ based on a single heating rate and with $R^2 > 0.974$ based on all heating rates. In addition, the predicted derivative thermogravimetry (DTG) profiles were integrated to obtain the TGA profiles, and the results agreed very well with the experimental data ($R^2 > 0.99$).

Keywords: Tobacco, Pyrolysis, Extra-Trees(ET), Machine learning

Nomenclature

TGA	Thermogravimetric analysis
DTG	Derivative thermogravimetry
ANN	Artificial neural networks
RF	Random forest
EC	Empirical correlations
ET	Extremely randomized trees
n_estimators	The number of trees in the forest
max_depth	TGA The maximum depth of each tree
TS	Total sugar
TA	Total alkaloids
RS	Reducing sugar
CL	Chlorine

K	Potassium
N	Nitrogen
CH	Chlorogenic acid
SC	Scopoletin
RU	Rutin
NE	Neochlorogenic acid
CR	Cryptochlorogenic acid
ST	Starch
PH	pH value
HR	Heating rate
T	Temperature
PCC	Pearson correlation coefficients
Mean	The average value
STD	The standard deviation
Min	The minimum value
Max	The maximum value
NRMSE	Normalized root mean square error
RMSE	Root mean square error

1

2 **1. Introduction**

3 It is a common sense that tobacco use is not good for health [1]. However,
 4 there is still a great deal of tobacco production, approximately 6.68 million
 5 metric tons in 2019 worldwide [2]. Thus, there is an urgent need to find
 6 other uses for tobacco instead of producing cigarettes. The compositions of

7 tobacco are cellulose, hemicellulose, lignin, and volatile component [3, 4, 5],
8 which are different from the lignocellulosic biomass that is mainly composed
9 of cellulose, hemicellulose, and lignin [6]. Although tobacco, as a special
10 kind of biomass, has different consistent of lignocellulosic biomass. It could
11 also be utilized by thermo-chemical and bio-chemical routes [7]. Among the
12 commonly-used thermo-chemical technologies, such as gasification, combus-
13 tion, and carbonation [3, 8, 9], pyrolysis is always the initial and fundamental
14 process. Therefore, it is of great significance to comprehensively study the
15 pyrolysis behaviour of tobacco and then accurately model this process for
16 developing healthy and efficient tobacco utilization technologies.

17 There have been some prior studies reported for tobacco pyrolysis. Ex-
18 perimental methods, such as gas chromatography/mass spectrometry [3, 10],
19 thermogravimetric analyzer coupled with Fourier transform infrared spec-
20 trometry [10, 11, 12], thermogravimetric mass spectrometry [3, 10] and macro-
21 thermogravimetric [13], were often used to study the pyrolysis behavior of
22 tobacco, including the effect of oxygen on in-situ evolution of chemical struc-
23 tures during tobacco pyrolysis [3], the pyrolysis characteristics and major py-
24 rolysis products [10], the effect of different tobacco particle sizes on pyrolysis
25 [14], and the formation mechanism of the volatile products [13]. Contrarily,
26 there have been only a few studies on establishing tobacco pyrolysis models.
27 Chen et al. [8] proposed a mathematical model for smouldering cigarettes,
28 including the evaporation of water, the pyrolysis of virgin tobacco and char-
29 combustion processes. However, the model only worked for the selected to-
30 bacco type, and its capability for various tobacco types was unclear. Encinar
31 et al. [15] established a kinetic model for pyrolysis (including maize, sun-

32 flowers, grape, and tobacco) from experiments. In this kinetic model, the
33 generation of volatile gases during the pyrolysis was assumed to be domi-
34 nated by a series of independent first-order parallel reactions. Gao et al. [16]
35 calculated the best linearity of different reaction models with the data from
36 experiments as a reference. However, tobacco pyrolysis involves very compli-
37 cated heterogeneous reactions, which relates to the tobacco type and heating
38 conditions. Therefore, it is hard to use one reaction or simple linear reactions
39 to represent the whole process. In summary, there is a lack of an accurate
40 kinetic model that considers both the effects of the chemical composition of
41 tobacco and the heating conditions.

42 The barrier to establishing such a kinetic model is the strong, complicated,
43 and nonlinear relationship between the pyrolysis behavior and the complex
44 chemical composition/heating conditions. Machine learning approaches have
45 been proven to handle complex nonlinear issues competently, as demon-
46 strated in many previous studies of solid fuel thermo-chemical conversion
47 modelling. Abbas et al. [17] proposed an approach based on Artificial Neu-
48 ral Networks (ANN) to predict the volatiles released from coal and biomass.
49 Yildiz et al. [18] studied the co-combustion of the hazelnut husk–lignite
50 coal blends of various compositions by using ANN. Sunphorka et al. [19]
51 modelled the biomass devolatilization process based on the three major com-
52 ponents using ANN, and the obtained model showed better performance than
53 conventional analytical method. Except for the ANN, ensemble learning ap-
54 proaches, such as Random Forest (RF), have also been successfully applied
55 in volatile gases yields and compositions prediction [20, 21], bio-oil yield pre-
56 diction [22], and biomass devolatilization kinetics prediction [23]. Xing et

57 al. [23] compared Empirical Correlations (EC), ANN, and RF in predicting
58 the kinetic parameters of a single-step kinetic model for biomass pyrolysis,
59 and found that RF gave the best performance. Wei et al. [24] compared RF
60 and Extremely Randomized Trees (ET) in the modelling of biomass and coal
61 co-pyrolysis.

62 Based on the above backgrounds, the objective of the present study is
63 to establish a general tobacco pyrolysis model based on the chemical con-
64 stituents and the heating conditions using an advanced machine learning
65 approach, ET. Specifically, in the experimental part, the chemical analysis
66 of 49 kinds of tobaccos were conducted, followed by the thermogravimet-
67 ric analysis (TGA) of tobaccos under a wide range of heating rates. The
68 derivative thermogravimetry (DTG) results were calculated by differentiat-
69 ing TGA. Subsequently, the chemical analysis results, the heating conditions,
70 and the DTG results are collected to construct the database. Afterwards,
71 the database is divided into train and test data-sets with a stratified sam-
72 pling method for the model development and evaluation. General pyrolysis
73 models were then developed based on the train data-set using ET. The fea-
74 ture selection and hyper-parameters optimization were conducted to get a
75 better generality and avoid over-fitting. The model performance was further
76 evaluated on the test data-set by comparing with the experimental data.

77 This work has the following novelty. Firstly, this is the first model directly
78 established the complex relation between tobacco samples' chemical compo-
79 sition and the pyrolysis behaviors. Chemical analysis is used to describe the
80 tobacco sample rather than proximate analysis and ultimate analysis. This
81 is because that there is a strong relation between chemical analysis with the

82 quality of tobacco, which is very important for the tobacco industry. Sec-
83 ondly, the differential target was used to improve the accuracy, as a small
84 error in TGA will be magnified to DTG by differentiating with a small step.
85 Thirdly, the present work established a database with a large amount range
86 of tobacco samples and heating rates, including 49 different kinds of tobacco
87 samples, and each of them are under a wide range of heating rate (10, 50,
88 100, 150, 200, 250, 300, 350, and 400 K/min). This database contributes to
89 the high generalization performance of model.

90 **2. Materials and methods**

91 In this section, the chemical analysis methods, the TGA experiment, and
92 the machine learning approaches are introduced in turn.

93 *2.1. Experimental methods*

94 In total 49 tobacco samples were studied in the experiments. The chem-
95 ical analysis information, including total sugar (TS), total alkaloids (TA),
96 reducing sugar (RS), chlorine (Cl), potassium (K), nitrogen (N), chlorogenic
97 acid (CH), scopoletin (SC), rutin (RU), neochlorogenic acid (NE), cryp-
98 tochlorogenic acid (CR), starch (ST), and pH value (PH), were measured.
99 The contents of TS, TA, RS, Cl, K, N, and ST were determined by a con-
100 tinuous flow analyzer (Alliance-Futura). The contents of CH, SC, RU, NE,
101 and CR were measured using an Agilent Technologies 1260 High-Performance
102 Liquid Chromatography system equipped with a DAD detector, in which the
103 separation process was performed on a Symmetry C18 column and the detec-
104 tion wavelength was 340 nm. PH was measured by the Mettler-Toledo Seven
105 Compact pH meter. The distribution densities of the chemical analysis of all

the 49 tobacco samples are shown in Fig. 1. The horizontal axis represents the feature distribution interval, and the vertical axis represents the distribution density. The curve in each sub-figure shows the distribution of the chemical analysis value. TS, TA, RS, CL, K, N, CH, SC, RU, NE, CR, ST and PH are within the ranges of 19.22-34.10 %, 1.48-3.27 %, 17.18-30.82 %, 0.07-1.10 %, 1.36-2.89 %, 1.63-2.42 %, 7.75-17.77 %, 0.14-049 %, 4.62-14.53 %, 1.13-2.92 %, 1.81-4.29 %, 3.41-7.17 % and 4.71-5.25, respectively. Table 1 lists some statistical information of the features, including the average value (Mean), the standard deviation (STD), the minimum (Min) and the maximum (Max). Obviously, TS has the largest mean value while SC has the smallest. STD is a measure of the amount of variation of a set of values, which lower STD means the values are closer to the average value. Except for TS, RS, CH and RU, other features are closer to their average value. More details are presented in supplementary materials. In total, there are 13 features used to describe the chemical constituents information of tobaccos, which are different from the traditional features used in biomass pyrolysis modelling such as proximate analysis and ultimate analysis. Chemical analysis is used to describe the tobacco sample rather than proximate analysis and ultimate analysis. This is because that there is a strong relation between chemical analysis with the quality of tobacco, which is very important for the tobacco industry[25].

Same as our previous study [4], all thermogravimetric analysis (TGA) of tobacco samples were conducted by using a *Discovery* thermogravimetric analyzer produced by TA instruments, as shown in Fig. 2, in the following way. Firstly, the tobacco samples were grounded into powders, and the particles

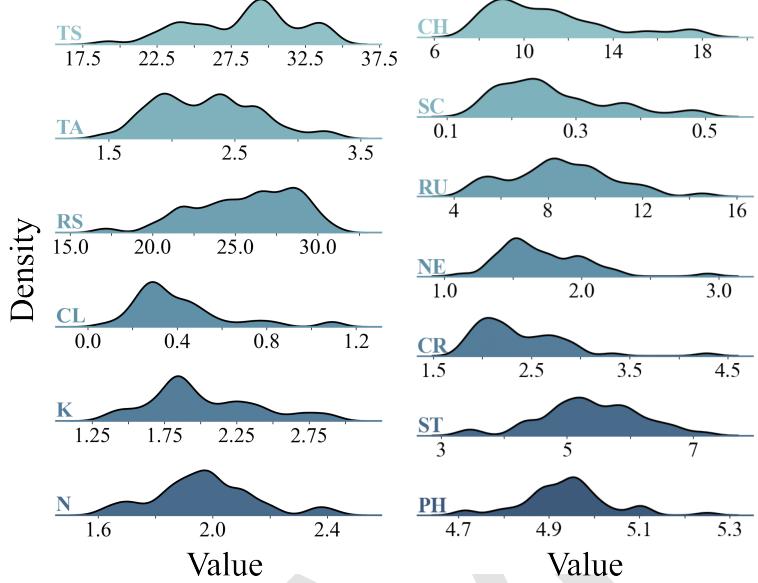


Figure 1: Distribution of tobacco chemical analysis. On the left side: TS, VA, RS, CL, K, N; on the right side: CH, SC, RU, NCH, ICH, ST, PH.

131 that were trapped on a 60-mesh sieve after passing through a 40-mesh sieve
 132 were collected in a sealed valve bag for the subsequent measurements. The
 133 size of tobacco particles is within the range of 0.250-0.425 mm. Secondly,
 134 8 mg samples were dehydrated by heating to 373 K at a heating rate of 10
 135 K/min from room temperature and holding for 30 minutes. The weight of
 136 samples after dehydration was set as the initial mass of the TGA measure-
 137 ment. Thirdly, the dehydrated samples were then devolatilized at different
 138 heating rates (10, 50, 100, 150, 200, 250, 300, 350, and 400 K/min) from
 139 373 K to 1173 K, and their mass loss with the increasing temperature was
 140 recorded. Among all tests, the flow rates of the carrier gas (high-purity N₂)
 141 and protective gas (high-purity N₂) were set at 50 mL/min and 30 mL/min,
 142 respectively. Each test has been repeated three times to ensure the repeata-

Table 1: The sample distribution of the test data-sets.

	Mean	STD	Min	Max
TS	28.31	3.66	19.22	34.10
TA	2.27	0.41	1.48	3.27
RS	25.67	3.10	17.18	30.82
CL	0.41	0.22	0.07	1.10
K	2.02	0.38	1.36	2.89
N	1.97	0.17	1.63	2.42
CH	11.04	2.65	7.75	17.77
SC	0.26	0.09	0.14	0.49
RU	8.54	2.20	4.62	14.53
NE	1.71	0.32	1.13	2.92
CR	2.37	0.46	1.81	4.29
ST	5.38	0.79	3.41	7.17
PH	4.94	0.10	4.71	5.25

143 bility of the results. Note that keeping the same sampling points during the
 144 experiment at both low and high heating rates is hard. Therefore, in or-
 145 der to avoid the uncertainty caused by different sampling points during the
 146 experiment, all the TGA results are interpolated into 500 points, and then
 147 the corresponding DTG results are calculated by differentiating TGA profiles
 148 based on particle temperature with Eq. (1),

$$DTG = \frac{TGA^{(i)} - TGA^{(i-1)}}{T^{(i)} - T^{(i-1)}}, \quad (1)$$

149 where $TGA^{(i)}$ and $T^{(i)}$ represent the mass loss of tobacco and local particle
 150 temperature at point i , respectively. $TG^{(i-1)}$ and $T^{(i-1)}$ represent the mass
 151 loss of tobacco and local particle temperature at point $(i - 1)$, respectively.

152 Worth noting that the present work is remarkably different from the pre-



Figure 2: Discovery thermogravimetric analyzer.

153 vious studies of pyrolysis modelling directly based on TGA profiles. On the
154 one side, empirically, building model follows the experimental way, which
155 first measures TGA by thermogravimetric analyzer then calculates DTG by
156 differentiating TGA using Eq. 1. On the other side, TGA is easier to model
157 than DTG owing to its simpler curve characteristics. However, there could
158 be a risk that some information would be lost if a big temperature step is
159 chosen in the calculation from TGA to DTG. Therefore, a small tempera-
160 ture step is usually preferred. A small difference in TGA will be magnified
161 to DTG by differentiating with a small step. Alternatively, owing to the
162 value of DTG being very small (usually in the range from 0 to -1), when
163 integrating DTG to TGA, the error will be reduced. Moreover, as mentioned
164 before, pyrolysis is the initial and fundamental process of the thermochemi-
165 cal conversion technologies of tobacco, in which the pyrolysis model provides
166 two-way coupling terms between the gas and particle phases in the numerical
167 simulations [26, 27]. Based on all discussion above, for avoiding the possible
168 uncertainties caused by differentiating TGA results and the convenience of

169 the model implementation in numerical simulations, the present work aims to
170 establish general tobacco pyrolysis models using the DTG results considering
171 the complex chemical constituents and heating conditions.

172 *2.2. Machine learning approaches*

173 The Extremely Randomized Trees (or called as Extra-Trees, ET) method
174 [28] is used to describe the complex nonlinear correlations between the pyrol-
175 ysis behavior and chemical constituents/heating conditions. The schematic
176 topological of the ET algorithmic can be seen in Fig. 3. This method has
177 been successfully applied for coal and biomass co-pyrolysis modelling in our
178 previous study [24], and therefore is briefly introduced here. Firstly, the
179 whole samples in the train data-set are inputted as learning samples rather
180 than using a bootstrap replica like the traditional methods (such as the ran-
181 dom forest). Then the trees will grow by splitting into two nodes as shown
182 in Fig. 3 with a threshold. In the traditional methods, the thresholds are
183 calculated, and then the best one is chosen. For ET, the splitting of each tree
184 is based on a random threshold. When the samples of the node are less than
185 the minimum sample size, the splitting stops. Finally, the predictions of the
186 trees are aggregated by an arithmetic average to yield the final prediction.

187 **3. Results and discussions**

188 In this section, firstly, the training processes, including feature selection
189 and hyper-parameters optimization, are introduced. Then the performance
190 of models on the train and test data-sets, and the further application from
191 the predicted DTG profiles to TGA profiles, are discussed in turn.

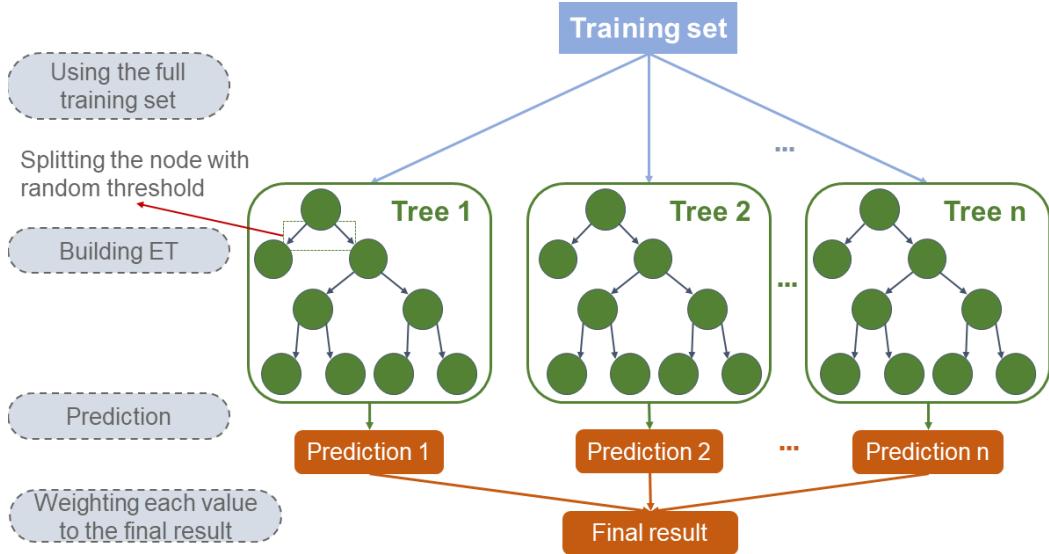


Figure 3: The schematic topological of the ET algorithmic [24].

192 3.1. Database information

193 In order to avoid the risk of introducing significant sampling bias when
 194 creating the test data-set by a purely random sampling method, herein, the
 195 stratified sampling method based on the difference of tobaccos is used[29]. In
 196 order to evaluate the difference of tobaccos, a parameter named Normalized
 197 Root Mean Square Error (NRMSE) is defined, which can be calculated by
 198 Eqs. (2) and (3),

$$RMSE_n = \left(\frac{1}{N} \sum_{i=1}^N \left(y_{1,\text{exp}}^{(i)} - y_{n,\text{exp}}^{(i)} \right)^2 \right)^{\frac{1}{2}}, \quad (2)$$

199

$$NRMSE = \frac{RMSE_n - RMSE_{n,\min}}{RMSE_{n,\max} - RMSE_{n,\min}}, \quad (3)$$

200 where $y_{1,\text{exp}}^{(i)}$ is the experimental value of tobacco 1 (base sample), and $y_{n,\text{exp}}^{(i)}$
 201 is the experimental value of tobacco n from 2 to 49. The base sample used in

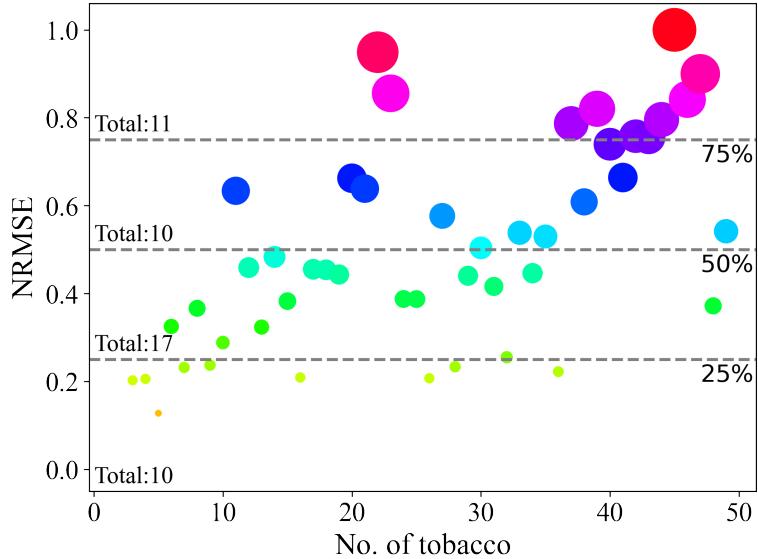


Figure 4: Results of the stratified sampling method, in which the size and color of the scatter vary with the values of NRMSE for a clearer observation.

the calculation is the experimental result of tobacco No.1 under 400 K/min, and the differences between the other 48 samples under 400 K/min with tobacco No.1 are shown in Fig. 4. The reason for choosing 400 K/min is that with the increase of heating rate, pyrolysis becomes more complicated and different [4], and the highest heating rate in the present work is 400 K/min. The different sizes and colors are used for a clearer observation. The size and color of the circles are determined by the NRMSE. The larger size stands for the larger NRMSE. Three grey dotted lines are used to divide the area (NRMSE from 0 to 1) into four equal parts, and the total sample numbers of each part are also shown. All the samples could be categorized into four types using the NRMSE values of 25%, 50%, and 75%. With this treatment, samples in each type are then randomly selected for training at a ratio of

214 70%, which could reflect the differences of tobacco pyrolysis behaviors and
215 further improve the model performances. The detailed information is shown
in Table 2.

Table 2: The sample distribution of the test data-sets.

<i>NRMSE</i>	Test	Total
0-25%	2,10,36	3
25-50%	13,15,19,24,31	5
50-75%	12,20,27	3
75-100%	22,37,40,42	4

216

217 3.2. Training process

218 After the train and test data-sets are chosen, the feature selection and
219 hyper-parameters optimization are conducted to train and optimize the model.
220 Worth noting that the present study established models not only based on
221 different heating rates (from 10 K/min to 400 K/min) but also on all heat-
222 ing rates. For convenience, the model built on the individual heating rate
223 is named with the heating rate, and the general model built on all heating
224 rates is named “all”.

225 3.2.1. Feature selection

226 As mentioned before, there are 13 features to describe the chemical com-
227 ponents of tobacco, as called the features of fuel properties, and 2 features,
228 temperature (T) and heating rate (HR), to describe the heating conditions.
229 In order to take a deep sight into the correlations between features of to-
230 baccos, the Pearson correlation coefficients (PCC) [30, 31] are determined

231 to find the relevance between them. The number in each grid is the PCC
 232 calculated by Eqs. (4) and (5),

$$\rho_{(X,Y)} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}, \quad (4)$$

233

$$\text{cov}(X, Y) = \mathbb{E} [(X - \mu_X)(Y - \mu_Y)], \quad (5)$$

234 where X and Y are the values of two different features, and $\text{cov}(X, Y)$ is the
 235 covariance between X and Y . σ_X , σ_Y , μ_X and μ_Y is the standard deviation
 236 of X , the standard deviation of Y , the mean of X and the mean of Y ,
 respectively.

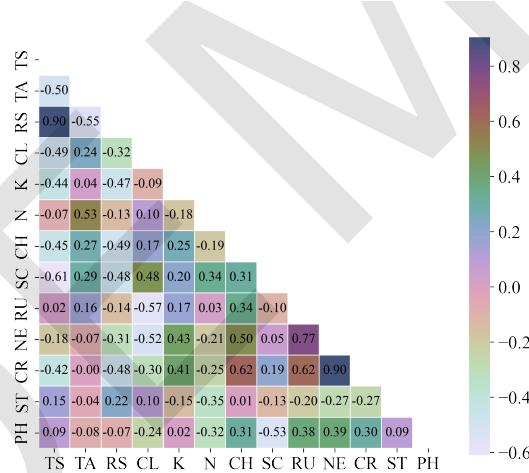


Figure 5: Pearson correlation coefficients between chemical constituents features of tobacco.

237

238 The PCC represents the strength and direction of linear relationships
 239 between pairs of features. According to Eqs. (4) and (5), The PCC of certain
 240 two features measure their linear correlation and has a value from -1 to
 241 1. Positive and negative values of PCC represent the positive and negative

correlations respectively, and the absolute values represent the strength of the correlation (the large values, the stronger correlation). For the training of machine learning models, it is recommended to select input variables that are not highly correlated to decrease the nonlinear degree of the system. For example, TS and RS are highly correlated since their PCC has a value of 0.93. So, the RS is deleted and the TS is kept. CR is also highly correlated with NE, CH, and RU, with PCC values equaling to 0.9, 0.9, 0.62, and 0.62, respectively. So, the NE, CH, and RU are deleted and CR is kept. Finally, the total number of features for the following machine learning modelling is 11, including 9 fuel properties features and 2 heating condition features.

3.2.2. Hyper-parameters optimization

Another important method to avoid over-fitting and reduce computational cost is hyper-parameters optimization. There are several hyper-parameters for ET in *scikit learn* [32], including *n_estimators* (the number of trees in the forest) and *max_depth* (the maximum depth of each tree). Herein, *n_estimators* and *max_depth* are determined by step-wise searching methods. The searching interval is 1 for both the two parameters, and the searching range is 1-150 and 1-30 for *n_estimators* and *max_depth*, respectively. The five-fold cross-validation method is used. This method randomly divides the train data-set into 5 sub-sets of equal size named folds. Each sub-set is used as a validation data-set to test the model, whereas the left 4 sub-sets are used for training [29]. Model performance on the validation data-set is used to determine the optimal hyper-parameters. For brevity, here only hyper-parameters optimization results of the model based on data

266 of tobacco pyrolysis under all heating rates are shown in Fig. 6, in which the
 267 arrows indicate the chosen parameters. It is found that with the increase
 268 of *n_estimators* and *max_depth*, the training and validation RMSEs sharply
 269 decrease and stabilize at a critical value, which means that the model per-
 270 formance would not be improved with the further increase of *n_estimators*
 271 and *max_depth*. Worth noting that there is a zig-zag pattern in the case of
 272 *n_estimators* as compared to *max_depth* optimization. This is because that
 273 the *n_estimators* has a relatively larger influence than that of *max_depth* in
 274 the Extra-tree algorithm. So, when both hyper-parameters begin to increase
 275 in the first place, n-estimation shows a zig-zag pattern and becomes smooth
 276 later. Therefore, the model with the critical hyper-parameters values (70 and
 277 14 for *n_estimators* and *max_depth*) is chosen and retained in the following
 278 discussions.

279 *3.3. Model performance*

280 In this section, the performances of models on both the train and test
 281 data-sets are presented and discussed. It's worth noting that all the tobacco
 282 on the test data-set is different from the train data-set, so they are totally
 283 unseen to the trained models. The *RMSE* and *R*² calculated by Eqs. (6), (7)
 284 and (8) are used to measure the performance of models,

$$RMSE = \left(\frac{1}{N} \sum_{i=1}^N \left(y_{\text{exp}}^{(i)} - y_{\text{pred}}^{(i)} \right)^2 \right)^{\frac{1}{2}}, \quad (6)$$

285

$$R^2 = 1 - \frac{\sum_{i=1}^N \left(y_{\text{pred}}^{(i)} - y_{\text{exp}}^{(i)} \right)^2}{\sum_{i=1}^N \left(y_{\text{exp}}^{(i)} - \hat{y}_{\text{exp}} \right)^2}, \quad (7)$$

286

$$\hat{y}_{\text{exp}} = \frac{1}{N} \sum_{i=1}^N y_{\text{exp}}^{(i)}, \quad (8)$$

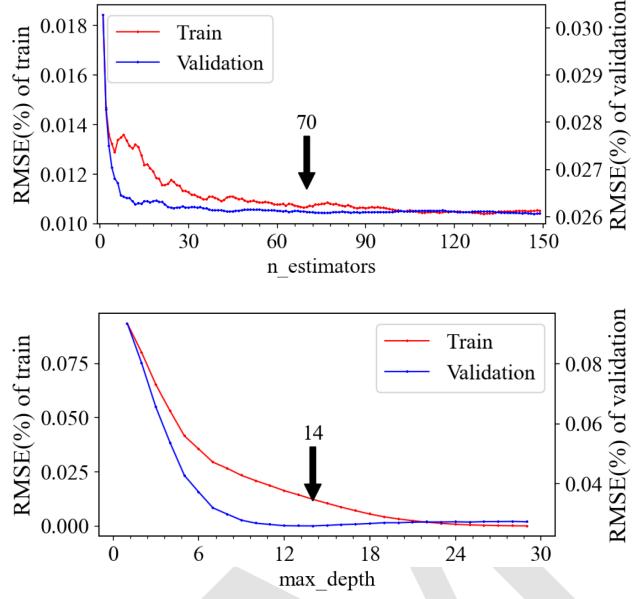


Figure 6: $N_{estimators}$ and max_depth optimization results of the trained model based on data from all heating rates.

where $y_{\text{exp}}^{(i)}$ is the experimental value, $y_{\text{pred}}^{(i)}$ is the predicted value and \hat{y}_{exp} is the average experimental value.

Figure 7 presents models performances on the train and test data-sets under heating rates of 10 K/min, 400 K/min, and all heating rates (including 10 K/min and 400 K/min). In which x axis stands for the true value while y axis stands for the predicted value from model, and different colors are used to present different heating rates (red, blue, and green for heating rate = 10, 400 K/min and all heating rates). In each subplot, the black dotted line is the best fitting line ($R^2 = 1$) and the dark line is the regression line of current data. Obviously, model based on 10 K/min shows the best performance, with R^2 on train test close to 1. While model based on 400 K/min shows

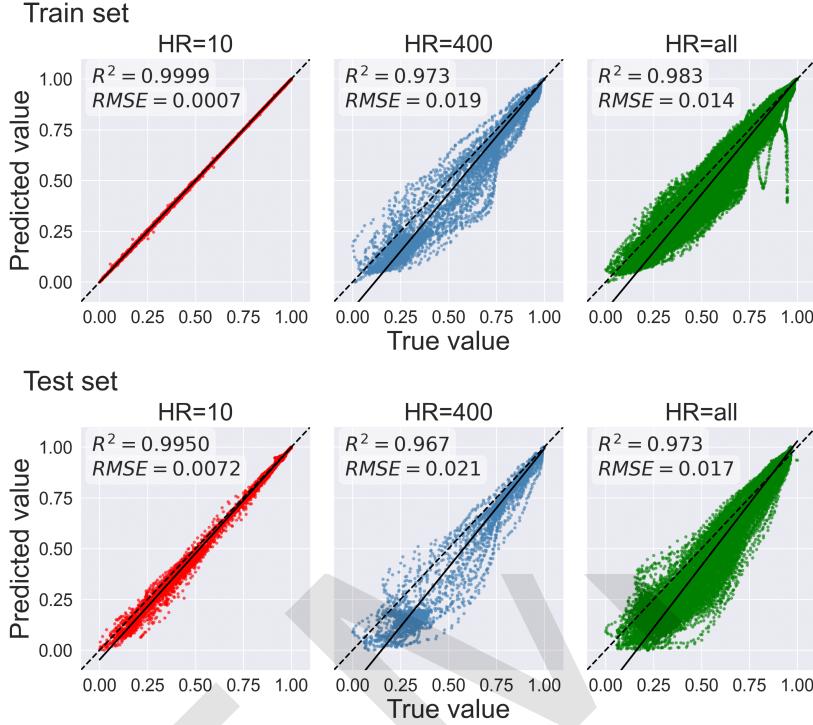


Figure 7: Model performances on the train and test data-sets under heating rates of 10 K/min, 400 K/min, and all heating rates.

the worst performance, which is even worse than model with all heating rate.
 More details are presented in Table 3, which lists the performance of models on train and test data-sets under different heating rates. It is found that 1)
 Models are not over-fitted. In all heating rates, the performances on the train data-set are comparable with those on the test data-set. The R^2 and $RMSE$ of predictions on the train data-set are only 0.009 larger and 0.0036 smaller than those on the test data-set at most. 2) Models perform well for both the train and test data-sets. The R^2 and $RMSE$ of the predictions of all the models on the test data-set are larger than 0.967 and smaller than 0.0209,

307 respectively. The best R^2 even equals to 0.995. 3) Inputting all heating rates
308 to train the model would lead to a good model generality. The performance
309 of “all” is not the best, but it’s better than some high heating rate models,
310 which indicates that the data-driven characteristics of the machine learning
311 approaches and more data inputting will improve the model’s generality. It’s
312 interesting to note that the predictions of the model based on 10 K/min have
313 the best performance on the test data-set with a R^2 value of 0.995, which is
314 much better than those of other models. With the increase of heating rate,
315 the performance of the model decreases. For the train data-set, both R^2 and
316 $RMSE$ decrease with the increase of heating rate, especially from 10 to 200
317 K/min. The test data-set shows the same trend. This could be attributed
318 to the fact that with the temperature increases, pyrolysis process becomes
319 more complicated with an aggregate phenomenon of the sub-peaks, which
320 was discovered in the work of Mu et al. [4]. According to Mu et al. [4],
321 the aggregate phenomenon under high heating rate is caused by the peak
322 positions of volatile components movement. Because volatile components’
323 pyrolysis is an exothermic reaction, the rate of heat accumulation around
324 tobacco particles will increase accordingly with the increase of heating rate,
325 which promotes the endothermic reaction and inhibits the exothermic reac-
326 tion. In addition, the peaks in DTG curve are related to the pyrolysis process
327 with different compositions of tobacco. Models for DTG curves with high
328 accuracy can be further used in studying the pyrolysis process of different
329 compositions of tobacco.

330 For a clearer assessment of model performance, both the best model (heat-
331 ing rate = 10K/min), and the worst model (heating rate=400 K/min) are

Table 3: Model performances on the train and test data-sets under different heating rates.

Heating rate (K/min)	Train		Test	
	R^2	$RMSE$	R^2	$RMSE$
10	0.999	0.0004	0.997	0.0060
50	0.992	0.0009	0.986	0.0119
100	0.985	0.0124	0.981	0.0140
150	0.989	0.0109	0.980	0.0147
200	0.981	0.0147	0.976	0.0165
250	0.973	0.0177	0.969	0.0201
300	0.970	0.0194	0.967	0.0203
350	0.973	0.0187	0.972	0.0109
400	0.970	0.0201	0.968	0.0208
All	0.990	0.0108	0.976	0.0164

332 further analyzed. Figures 8 and 9 show the comparisons between the DTG
 333 profiles predicted with the models and measured in the experiments for var-
 334 ious tobacco samples from the test data-set. The red line and blue scatters
 335 stand for the predicted results and the experimental data, respectively. The
 336 No., R^2 and $RMSE$ are also shown in each subplot. From Figure 8, the
 337 predicted values of model based on 10 K/min show great agreement with
 338 experimental values, with all R^2 s larger than 0.9909. It can be concluded
 339 that model on 10 K/min is at very high accuracy. Compared with model on
 340 10 K/min, there is a slightly decrease of accuracy for model based on 400
 341 K/min. Figure 9 demonstrates that predictions are not as good as those of
 342 former model, but generally in good agreement with the experimental data.
 343 The slight deviations are mainly concentrated in the peak area, where the
 344 temperature is within the range of 500 K to 700 K, and in other temperature

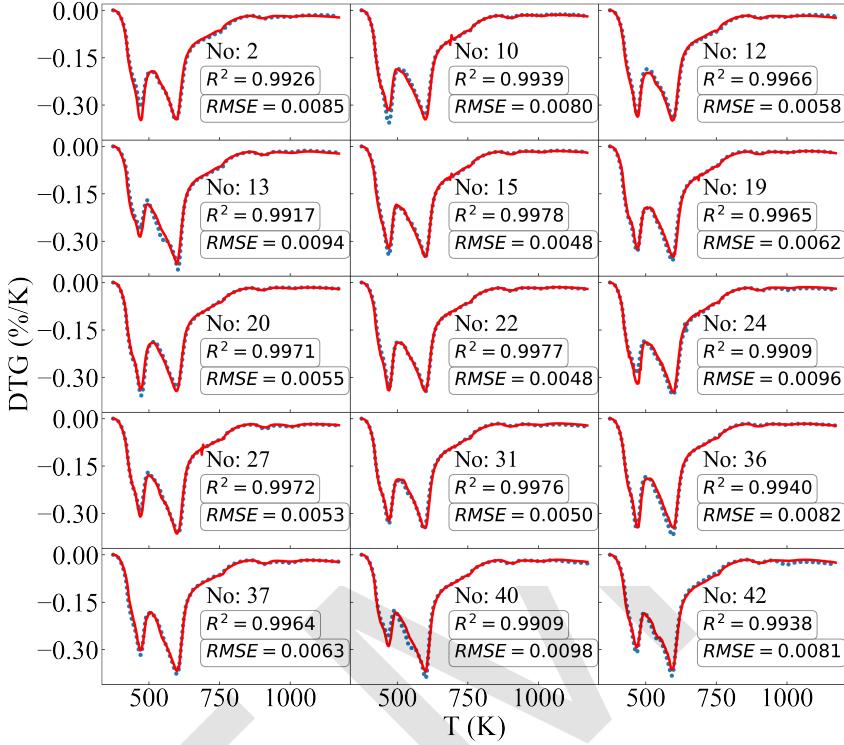


Figure 8: Comparisons between the DTG profiles predicted with the model and measured in the experiments for various tobacco samples at a heating rate of 10 K/min.

345 ranges, models performed well. In the peak area, the model can predict the
 346 peak position temperatures accurately. However, the predicted heights of the
 347 peaks are slightly lower than those in experiments.

348 To check how those deviations in the peak area affect the further cal-
 349 culations of TGA profiles, the predicted DTG at 400 K/min profiles are
 350 integrated based on temperature to obtain the TGA profiles. Figure 10
 351 shows the comparisons between TGA profiles obtained from the predicted
 352 DTG profiles and measured in the experiments for various tobacco samples.
 353 The figure legend is the same as that in Fig. 9. Overall, the predicted TGA

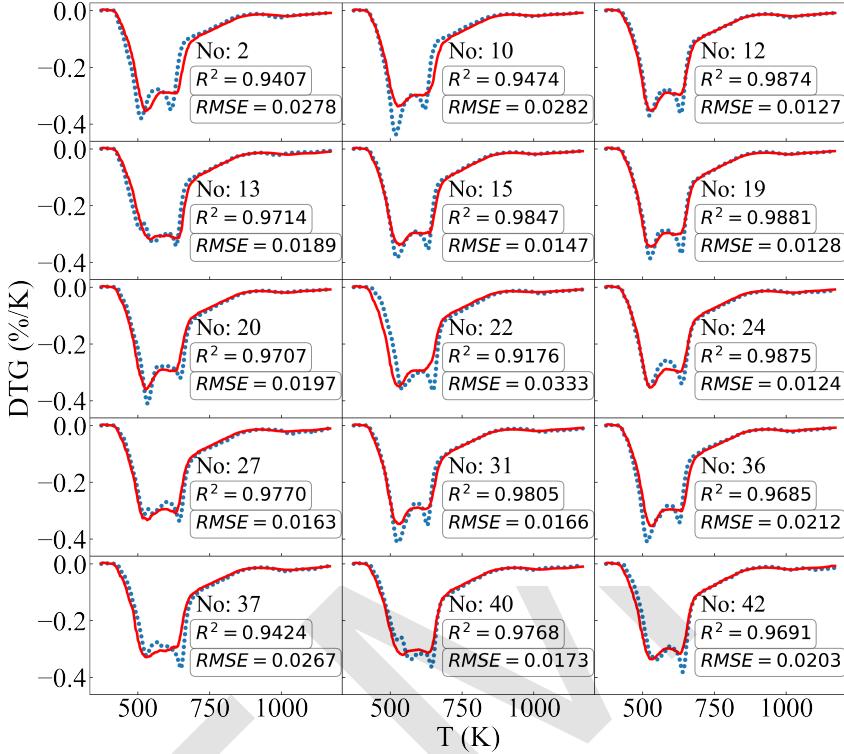


Figure 9: Comparisons between the DTG profiles predicted with the model and measured in the experiments for various tobacco samples at a heating rate of 400 K/min.

354 profiles are in a better agreement with the experimental data compared with
 355 those of the DTG profiles with all the R^2 larger than 0.99. This is because
 356 of the small DTG peak values and the slight deviations in the peak height
 357 prediction. In addition, the integration process also reduces the deviations
 358 statistics as it has been discussed in section 2.1. The peaks in the DTG
 359 profiles could be well reproduced with the developed models. Those peaks
 360 relate to the pyrolysis of different types of components, including the extrac-
 361 tive content, cellulose, hemicellulose and lignin, as described in our previous
 362 study [1]. Those components could also be used as the inputs in the model

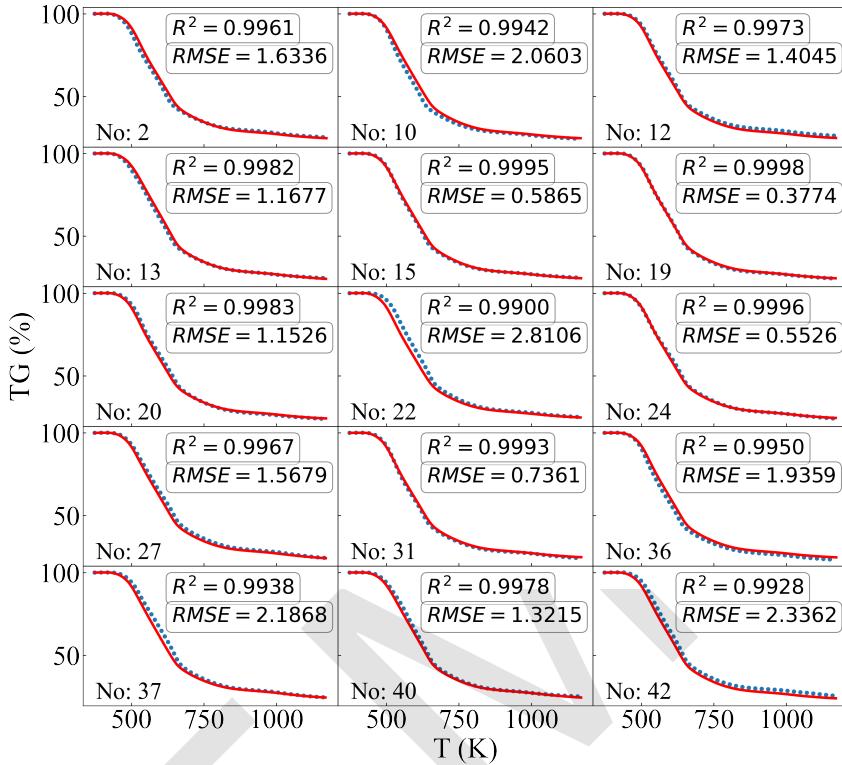


Figure 10: Comparisons between the TGA profiles integrated from the predicted DTG profiles and measured in the experiments for various tobacco samples at a heating rate of 400 K/min.

development, which will be further explored in the future study.

Present work currently has the following limitation. The heating rate range of the experimental equipment used in the present study (The discovery thermogravimetric analyzer produced by TA instruments) can't be over 400 K/min, so only the experiment under 400 K/min can be done. The model performance over 400 K/min is still unclear.

³⁶⁹ **4. Conclusions**

³⁷⁰ In the present study, a general tobacco pyrolysis model was developed
³⁷¹ based on the chemical constituents and heating conditions using the ET
³⁷² method. Specifically, thermogravimetric and chemical analysis of wide ranges
³⁷³ of tobacco types and heating conditions are first conducted by experiments
³⁷⁴ to construct a database for the model development. Subsequently, the con-
³⁷⁵ structed database has been divided into train and test data-sets for the model
³⁷⁶ development and evaluation. General devolatilization models were then de-
³⁷⁷ veloped based on the train data-set using ET. The performances of models
³⁷⁸ were further evaluated on the test data-set by comparing with the experimen-
³⁷⁹ tal data. The results showed that after feature selection based on Pearson
³⁸⁰ correlation coefficient and hyper-parameters optimization, the trained model
³⁸¹ could well reproduce the tobacco pyrolysis behavior on the test data with
³⁸² $R^2 > 0.967$ based on a single heating rate and with $R^2 > 0.974$ based on all
³⁸³ heating rates. In addition, the worst predicted DTG profiles were integrated
³⁸⁴ based on particle temperature to obtain the TGA profiles, and the results
³⁸⁵ showed a very well agreement with the experimental data ($R^2 > 0.99$). The
³⁸⁶ present work concluded that the data-driven model for tobacco pyrolysis can
³⁸⁷ achieve a very high accuracy and robustness, which provide a new way of
³⁸⁸ modeling pyrolysis. In addition, setting DTG (differenced TG) as the target
³⁸⁹ and then integrating it back to TG is a more accuracy for data-driven py-
³⁹⁰ rolysis models since a small difference in TGA will be magnified to DTG by
³⁹¹ differentiating with a small step.

³⁹² **Acknowledgments**

393 The authors are grateful for the support from National Natural and Sci-
394 ence Foundation of China (Grant No. 51925603) and the support from Zhe-
395 jiang University-Zhejiang China Tobacco Joint Laboratory Fund (Grant No.
396 K-20201802).

397 **References**

- 398 [1] D. Hammond, Health warning messages on tobacco products: a review,
399 *Tobacco control* 20 (5) (2011) 327–337.
- 400 [2] M. Shahbandeh, Topic: Tobacco industry.
401 URL <https://www.statista.com/topics/1593/tobacco/#dossierKeyfigures>
- 403 [3] Y. Peng, X. Hao, Q. Qi, X. Tang, Y. Mu, L. Zhang, F. Liao, H. Li,
404 Y. Shen, F. Du, et al., The effect of oxygen on in-situ evolution of
405 chemical structures during the autothermal process of tobacco, *Journal*
406 of Analytical and Applied Pyrolysis 159 (2021) 105321.
- 407 [4] Y. Mu, Y. Peng, X. Tang, J. Ren, J. Xing, K. Luo, J. Fan, K. Zhang, Ex-
408 perimental and kinetic studies on tobacco pyrolysis under a wide range
409 of heating rates, *ACS Omega* (2021) acsomega.1c06122.
- 410 [5] Y. J. Sung, Y. B. Seo, Thermogravimetric study on stem biomass of
411 *nicotiana tabacum*, *Thermochimica Acta* 486 (1-2) (2009) 1–4.
- 412 [6] S. Wang, G. Dai, H. Yang, Z. Luo, Lignocellulosic biomass pyrolysis
413 mechanism: a state-of-the-art review, *Progress in energy and combus-*
414 *tion science* 62 (2017) 33–86.

- 415 [7] S. Wang, G. Dai, H. Yang, Z. Luo, Lignocellulosic biomass pyrolysis
416 mechanism: A state-of-the-art review, *Progress in Energy and Combus-*
417 *tion Science* 62 (2017) 33–86.
- 418 [8] P. Chen, A mathematical model of cigarette smoldering process,
419 *Beiträge Zur Tabakforschung International/Contributions to Tobacco*
420 *Research* 20 (4) (2002) 265–271.
- 421 [9] J. Xing, Y. Bai, C. Zhao, Z. Gao, H. Wang, Numerical studies of coal
422 devolatilization characteristics with gas temperature fluctuation, *Energy*
423 & *Fuels* 32 (8) (2018) 8760–8767.
- 424 [10] B. Liu, Y.-M. Li, S.-B. Wu, Y.-H. Li, S.-S. Deng, Z.-L. Xia, Pyrolysis
425 characteristic of tobacco stem studied by py-gc/ms, tg-ftir, and tg-ms,
426 *BioResources* 8 (1) (2013) 220–230.
- 427 [11] R. R. Baker, L. J. Bishop, The pyrolysis of tobacco ingredients, *Journal*
428 *of analytical and applied pyrolysis* 71 (1) (2004) 223–311.
- 429 [12] E. Calabuig, N. Juárez-Serrano, A. Marcilla, Tg-ftir study of evolved gas
430 in the decomposition of different types of tobacco. effect of the addition
431 of sba-15, *Thermochimica Acta* 671 (2019) 209–219.
- 432 [13] F. Barontini, A. Tugnoli, V. Cozzani, J. Tetteh, M. Jarriault, I. Zinovik,
433 Volatile products formed in the thermal decomposition of a tobacco
434 substrate, *Industrial & Engineering Chemistry Research* 52 (42) (2013)
435 14984–14997.
- 436 [14] G. Guo, C. Liu, Y. Wang, S. Xie, K. Zhang, L. Chen, W. Zhu, M. Ding,
437 Comparative investigation on thermal degradation of flue-cured tobacco

- 438 with different particle sizes by a macro-thermogravimetric analyzer and
439 their apparent kinetics based on distributed activation energy model,
440 *Journal of Thermal Analysis and Calorimetry* 138 (5) (2019) 3375–3388.
- 441 [15] J. M. Encinar, F. J. Beltrán, J. F. González, M. J. Moreno, Pyrolysis
442 of maize, sunflower, grape and tobacco residues, *Journal of Chemical
443 Technology & Biotechnology: International Research in Process, Envi-
444 ronmental AND Clean Technology* 70 (4) (1997) 400–410.
- 445 [16] W. Gao, K. Chen, Z. Xiang, F. Yang, J. Zeng, J. Li, R. Yang, G. Rao,
446 H. Tao, Kinetic study on pyrolysis of tobacco residues from the cigarette
447 industry, *Industrial Crops and Products* 44 (2013) 152–157.
- 448 [17] T. Abbas, M. Awais, F. Lockwood, An artificial intelligence treatment
449 of devolatilization for pulverized coal and biomass in co-fired flames,
450 *Combustion and flame* 132 (3) (2003) 305–318.
- 451 [18] Z. Yıldız, H. Uzun, S. Ceylan, Y. Topcu, Application of artificial neural
452 networks to co-combustion of hazelnut husk–lignite coal blends, *Biore-
453 source technology* 200 (2016) 42–47.
- 454 [19] S. Sunphorka, B. Chalermsinsuwan, P. Piemsomboon, Artificial neu-
455 ral network model for the prediction of kinetic parameters of biomass
456 pyrolysis from its constituents, *Fuel* 193 (2017) 142–158.
- 457 [20] Q. Tang, Y. Chen, H. Yang, M. Liu, H. Xiao, S. Wang, H. Chen, S. R.
458 Naqvi, Machine learning prediction of pyrolytic gas yield and compo-
459 sitions with feature reduction methods: Effects of pyrolysis conditions
460 and biomass characteristics, *Bioresource Technology* 339 (2021) 125581.

- 461 [21] J. Xing, K. Luo, H. Wang, T. Jin, J. Fan, Novel sensitivity study for
462 biomass directional devolatilization by random forest models, Energy &
463 Fuels 34 (7) (2020) 8414–8423.
- 464 [22] Z. Ullah, S. R. Naqvi, W. Farooq, H. Yang, S. Wang, D.-V. N. Vo,
465 et al., A comparative study of machine learning methods for bio-oil
466 yield prediction—a genetic algorithm-based features selection, Biore-
467 source Technology 335 (2021) 125292.
- 468 [23] J. Xing, H. Wang, K. Luo, S. Wang, Y. Bai, J. Fan, Predictive single-step
469 kinetic model of biomass devolatilization for cfd applications: A com-
470 parison study of empirical correlations (ec), artificial neural networks
471 (ann) and random forest (rf), Renewable Energy 136 (2019) 104–114.
- 472 [24] H. Wei, K. Luo, J. Xing, J. Fan, Predicting co-pyrolysis of coal and
473 biomass using machine learning approaches, Fuel 310 (2022) 122248.
- 474 [25] H. Wei, Y. Peng, H. Huang, J. Fan, J. Xing, K. Luo, J. Fan,
475 L. Dai, Toba-cpd: An extended chemical percolation devolatilization
476 model for tobacco pyrolysis, ACS Omega 7 (41) (2022) 36776–36785.
477 arXiv:<https://doi.org/10.1021/acsomega.2c05098>, doi:10.1021/
478 acsomega.2c05098.
479 URL <https://doi.org/10.1021/acsomega.2c05098>
- 480 [26] J. Xing, K. Luo, H. Wang, Z. Gao, T. Jin, J. Fan, Comparative study
481 on different treatments of coal devolatilization for pulverized coal com-
482 bustion simulation, Energy & Fuels 34 (3) (2020) 3816–3827.

- 483 [27] C. Zhao, K. Luo, R. Cai, J. Xing, Z. Gao, J. Fan, Large eddy simulations
484 and analysis of no emission characteristics in a laboratory pulverized coal
485 flame, Fuel 279 (2020) 118316.
- 486 [28] P. Geurts, D. Ernst, L. Wehenkel, Extremely randomized trees, Machine
487 learning 63 (1) (2006) 3–42.
- 488 [29] A. Géron, Hands-on machine learning with Scikit-Learn, Keras, and
489 TensorFlow: Concepts, tools, and techniques to build intelligent sys-
490 tems, O'Reilly Media, 2019.
- 491 [30] Pearson correlation coefficient (2021).
- 492 [31] K. Pearson, Mathematical contributions to the theory of evolution.—on
493 a form of spurious correlation which may arise when indices are used in
494 the measurement of organs, Proceedings of the royal society of london
495 60 (359-367) (1897) 489–498.
- 496 [32] sklearn.ensemble.ExtraTreesRegressor.
497 URL <https://scikit-learn/stable/modules/generated/sklearn.ensemble.ExtraTreesRegressor.html>