



Kent Academic Repository

Sun, Caiying, Wang, Lijuan, Yan, Yong, Zhang, Wenbiao and Shao, Ding (2021) *A Novel Heterogeneous Ensemble Approach to Variable Selection For Gas-Liquid Two-Phase CO₂ Flow Metering*. *International Journal of Greenhouse Gas Control*, 110 . ISSN 1750-5836.

Downloaded from

<https://kar.kent.ac.uk/89452/> The University of Kent's Academic Repository KAR

The version of record is available from

<https://doi.org/10.1016/j.ijggc.2021.103418>

This document version

Author's Accepted Manuscript

DOI for this version

Licence for this version

CC BY-NC-ND (Attribution-NonCommercial-NoDerivatives)

Additional information

Versions of research works

Versions of Record

If this version is the version of record, it is the same as the published version available on the publisher's web site. Cite as the published version.

Author Accepted Manuscripts

If this document is identified as the Author Accepted Manuscript it is the version after peer review but before type setting, copy editing or publisher branding. Cite as Surname, Initial. (Year) 'Title of article'. To be published in *Title of Journal*, Volume and issue numbers [peer-reviewed accepted version]. Available at: DOI or URL (Accessed: date).

Enquiries

If you have questions about this document contact ResearchSupport@kent.ac.uk. Please include the URL of the record in KAR. If you believe that your, or a third party's rights have been compromised through this document please see our [Take Down policy](https://www.kent.ac.uk/guides/kar-the-kent-academic-repository#policies) (available from <https://www.kent.ac.uk/guides/kar-the-kent-academic-repository#policies>).

A Novel Heterogeneous Ensemble Approach to Variable Selection For Gas-Liquid Two-Phase CO₂ Flow Metering

Caiying Sun^{1,3}, Lijuan Wang², Yong Yan^{2,*}, Wenbiao Zhang¹, Ding Shao¹

¹ School of Control and Computer Engineering, North China Electric Power University, Beijing
102206, P.R. China

² School of Engineering and Digital Arts, University of Kent, Canterbury, Kent CT2 7NT, U.K.

³ School of Information Engineering, Inner Mongolia University of Science & Technology, Baotou
014010, P.R. China

* Correspondent author: y.yan@kent.ac.uk

Abstract

Variable selection is an important preprocessing step in the development of effective data-driven models for CO₂ flow measurement in carbon capture and storage systems. In order to effectively quantify the importance of potential input variables to the desired output, ensemble learning is proposed and incorporated into variable selection methodology. This paper presents a tree-based heterogeneous ensemble approach to variable selection and its application to gas-liquid two-phase CO₂ flow measurement. The importance of each variable is determined through combining the importance scores from four tree-based algorithms, including decision tree regression, bootstrap aggregating of regression trees, gradient boosting decision tree and gradient boosting random forest. Then the backward elimination algorithm is applied to remove the relatively less important variables and hence a small set of input variables for data-driven models. The selection results demonstrate that the significant variables for CO₂ mass flow measurement include *apparent mass flow rate*, *time shift*, *differential pressure* and *pressure drop* while *observed density*, *density drop*, *observed flow velocity* and *outlet temperature* for prediction of gas volume fraction. To assess the validity of the selected variables, data-driven models based on gradient boosting random forest are developed. Results suggest that the relative error of the model output is mostly within 1% for CO₂ mass flowrate measurement and 5% for gas volume fraction prediction by taking the selected variables as model inputs.

Keywords: carbon capture and storage, gas-liquid two-phase CO₂, variable selection, heterogeneous ensemble approach, data-driven models

34 1. Introduction

35 With the rapid development of machine learning technology, variable selection becomes more and
36 more important in data analysis and data-driven modelling (Wang et al., 2013; Xin et al., 2012; Nan
37 et al., 2014; Zhang et al., 2018). A variety of variable selection methods have been developed over
38 the past few years (Tuv et al., 2009; Zhang et al., 2015; Zhang et al., 2017). However, there is no
39 common rule to determine which method is suitable for a particular application. It is normally
40 determined by balancing the computational cost and the accuracy of the output from the data-driven
41 model. Zhu et al. (2006) extended the application of ensemble learning methods (Mendes-Moreira
42 et al., 2012; AL-Qutami et al., 2018) from the prediction ensemble to the variable selection
43 ensemble. In general, the ensemble approach for variable selection can be classified into two
44 categories: homogeneous and heterogeneous approaches. (Zhu et al., 2011; Zhou., 2012; Li et al.,
45 2017). The homogeneous ensemble approach is to use the same selection method on different
46 datasets while the heterogeneous ensemble approach is to train different selection algorithms on the
47 same dataset.

48

49 CO₂ flow in carbon capture and storage (CCS) systems is of complex nature (Wang et al., 2018;
50 Zhang et al., 2018; Shao et al., 2020) and it is thus challenging to measure its dynamic
51 characteristics. To measure mass flowrate and gas volume fraction of multi-phase flow, data-driven
52 modelling has been considered as an efficient and cost-effective way (Yan et al., 2018).
53 Applications of Coriolis flowmeters to gas-liquid two-phase flow measurement have been attempted
54 by using prototype transmitters and investigating into the use of the internal parameters (Green et al.,
55 2008; Kunze et al., 2014; Li et al., 2018; Li et al., 2019). Coriolis flowmeters incorporating
56 data-driven modelling algorithms have demonstrated a potential for multiphase flow measurement
57 (Wang et al. 2017). One key feature of this approach is to minimise the hardware modification and
58 enable commercial Coriolis flowmeters to work under two-phase flow conditions by simply adding
59 a software module. In order to develop optimal data-driven models for Coriolis flowmeters under
60 two-phase flow conditions and quantify the parametric dependency among the input variables and
61 their significance to the desired outputs, Wang et al. (2017) compared three input variable selection
62 methods, partial mutual information, genetic algorithm-artificial neural network, and tree-based
63 iterative input selection. It is found that a single tree-based selection method can generate varying

64 results for different datasets in the variable selection process.

65

66 To improve the performance of the tree-based selection method, a heterogeneous ensemble approach
67 is introduced to the variable selection process in this paper. A total of four different tree-based
68 selection methods, including decision tree (DT) regression, bootstrap aggregating (Bagging) of
69 regression trees, gradient boosting decision tree (GBDT) and gradient boosting random forest (GBRF)
70 are implemented with the same dataset and fused importance score is calculated for each variable.
71 This paper aims to propose an approach to input variable selection and data-driven modelling
72 for two-phase flow measurement and test the developed models on the same type of
73 flowmeters and transmitters. Experimental tests were conducted with gas-liquid two-phase CO₂
74 flow. The validity of the selected variables is verified by assessing the performance of GBRF based
75 data-driven models.

76

77 **2. Methodology**

78 The structure of the ensemble variable selection method is shown in Fig. 1. The dataset is acquired
79 from multiple sensors including a Coriolis mass flowmeter, a DP transducer, two pressure
80 transducers and two temperature sensors. The first step is to generate variable selectors with
81 optimized parameters based on different tree-based algorithms. In this step, the importance scores of
82 variables are obtained from each selector. The second step is to combine the variable importance
83 derived from different variable selectors and then remove the less important variables through
84 backward elimination. Therefore, a set of variables which has significant effect on the mass flow rate
85 measurement and gas volume fraction (GVF) prediction of two-phase CO₂ flow is obtained,
86 respectively. The third step is to develop data-driven models based on the selected variables to
87 produce mass flow rate and GVF.

88

89

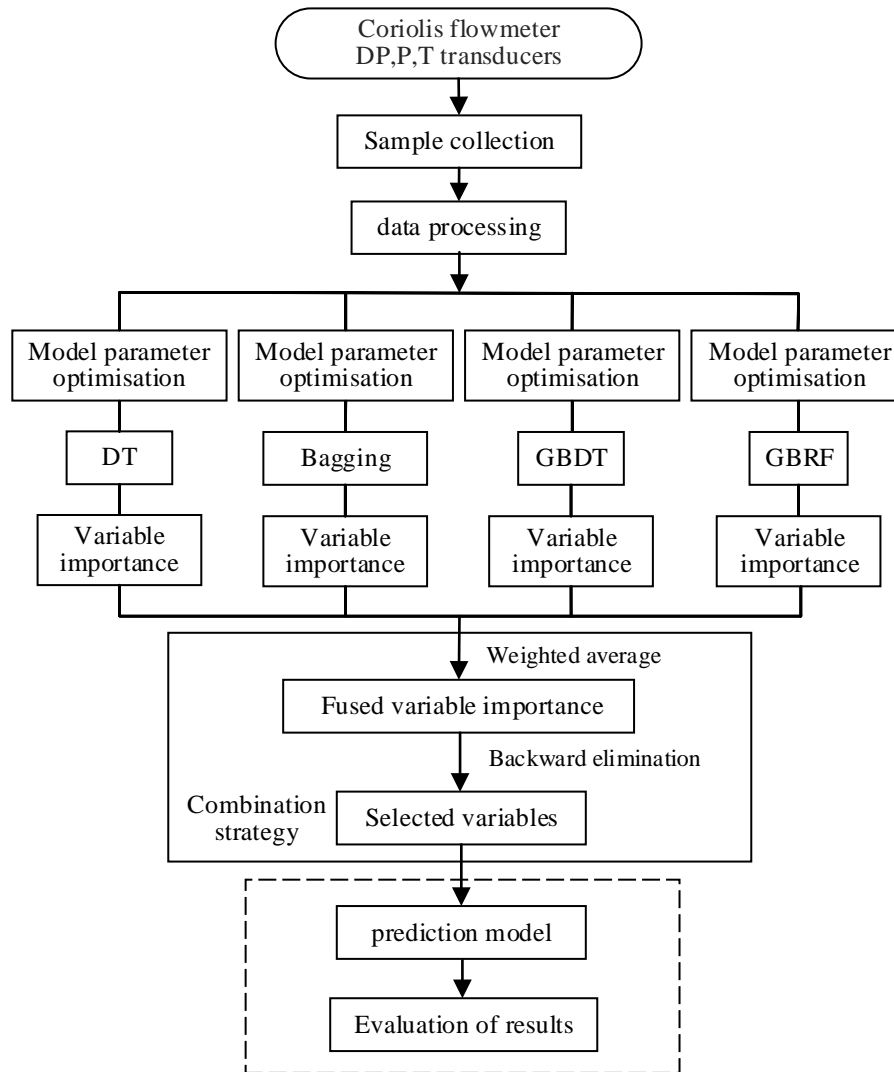


Fig. 1 Structure of the ensemble variable selection method

90

91

92

93 2.1 Variable selectors

94 Four tree-based models, including DT, Bagging, GBDT and GBRF, are taken as four individual
 95 selectors. These are commonly used tree based algorithms and are effective in the variable selection.

96 The importance score of each variable can be quantified by all tree-based algorithms, respectively. The

97 weighted average method is then used to combine the importance scores from different tree based

98 models. The application of ensemble learning is to improve the reliability of variable selection. As

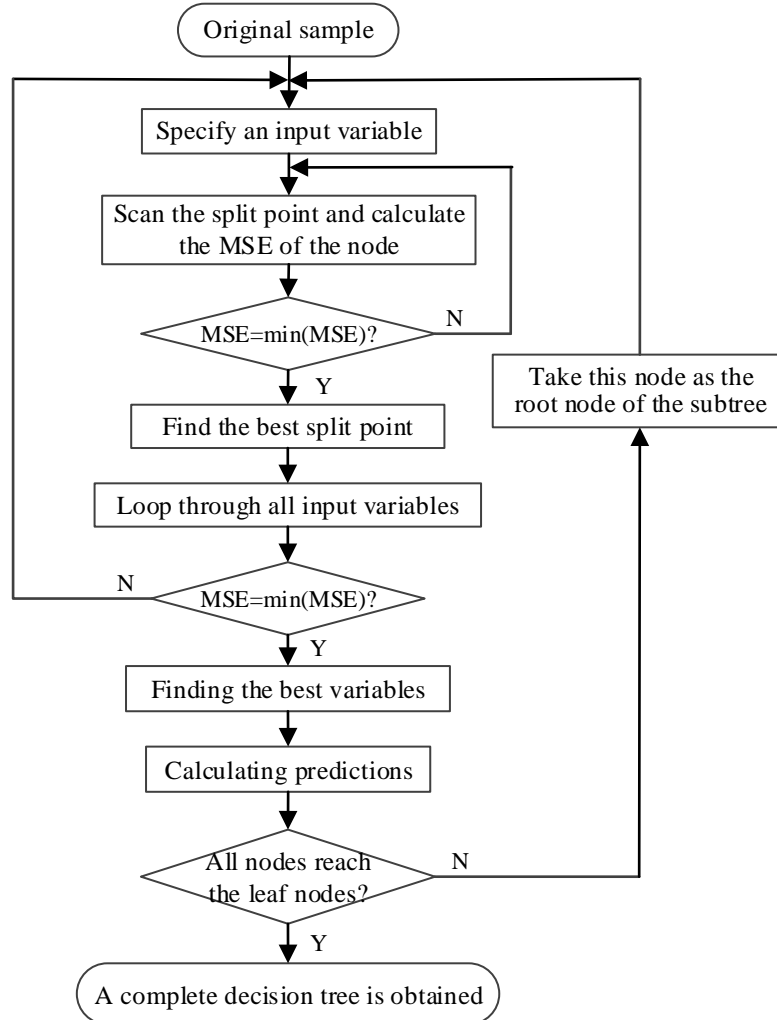
99 shown in Fig. 2, the DT algorithm (Zhou., 2012) traverses all the input variables in each iteration and

100 take the input variable that produces the minimum MSE (Mean Squared Error) value of the prediction

101 result as the split point of the node. The process continues recursively until the row arrives at a

102 terminal (leaf) node where a prediction value is assigned to the row. The value assigned to the terminal

103 node is the mean of the outcomes of all training observations that wound up in the leaf node. These
 104 observation results are the predicted values corresponding to the predicted target parameters. In this
 105 paper, the observation results are mass flow rate and GVF of gas-liquid two-phase CO₂ flow.
 106



107
 108 Fig. 2 Decision tree flowchart

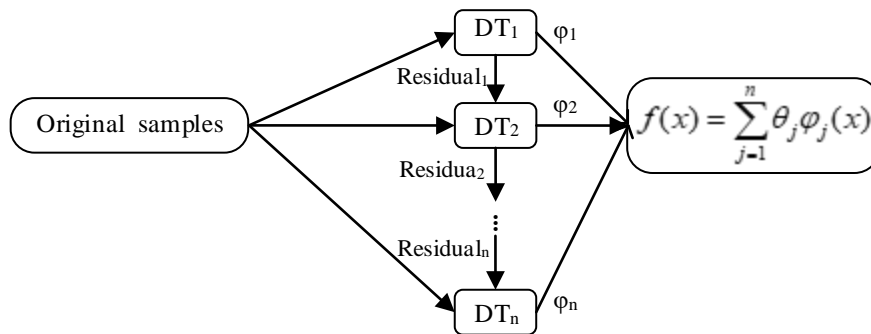
109 For a single decision tree, the more important the variable is, the earlier it is used. The variable
 110 importance is quantified as (Tuv et al., 2009):

111
$$VI(X_i, T) = \sum_{t \in T} \Delta I(X_i, t) \quad (1)$$

112 where $\Delta I_i(X_i, t)$ is the reduction in impurity due to an actual (or potential) split on the variable X_i at the
 113 node t of the optimally pruned tree T . Node impurity $I(t)$ in this paper is MSE with a node t . After
 114 normalizing the VI of each variable, the final importance score will be obtained.

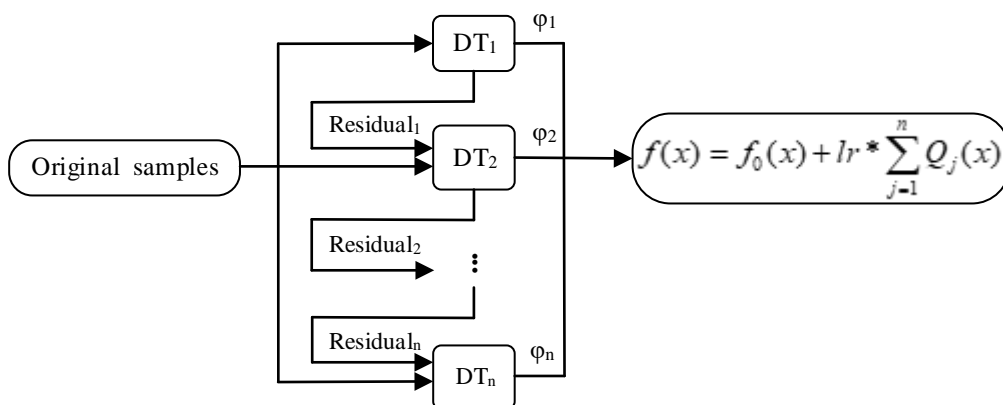
115
 116 Bagging algorithm (Zhang et al., 2015; Zhu et al., 2011) uses bootstrap sampling to obtain the data

117 subsets for training the base learners (i.e. decision trees). In Fig. 3, φ_i is the prediction result of the i^{th}
 118 decision tree for the test samples. The output $f(x)$ in regression is the average of the results from all
 119 trees. For tree-based ensembles, the importance score of a variable is the average value derived from
 120 all the trees.



121
 122 Fig. 3 Block diagram of Bagging algorithm

123
 124 The GBDT (Zhang et al., 2015; Tuv et al., 2009) algorithm is a decision tree ensemble learning
 125 algorithm based on the gradient boosting framework. It uses the original training samples for learning.
 126 The i^{th} residual in Fig. 4 is the difference between the predicted result of the i^{th} tree and the target value.
 127 It is taken as a new target for the next tree to achieve the concatenation of GBDT algorithm. The
 128 output $f(x)$ of the algorithm is obtained by summing the prediction results of individual trees. $f_0(x)$ is
 129 the initial value of the learner. lr is the learning rate, which is used to control the step size and ensure
 130 the convergence of the algorithm during the iteration.



131
 132 Fig. 4 Block diagram of GBDT algorithm

133

134 GBRF algorithm (Tuv et al., 2009) is a mixture of gradient boosting and random forest algorithms.
 135 When dividing the samples at each node of the tree in GBRF algorithm, it only takes max-features
 136 attributes at random rather than all attributes. The GBRF algorithm introduces random selection of
 137 input variables in the splitting process and thus the correlation between the single models is further
 138 reduced. The structure of the GBRF algorithm is the same as that of GBDT. By changing the value of
 139 different coefficients in the model, different ranking results of variables are obtained.

140

141 2.2 Combination strategy

142 Combination strategy which is used to fuse importance score from individual selectors usually
 143 includes averaging method, voting method and learning method. Stacking (Zhou., 2012; Breiman.,
 144 1996), as a typical learning method, is to train the first-level selectors using the original training
 145 dataset. The fused importance FI for a particular variable is the weighted average of importance scores
 146 from individual selectors and defined as:

$$147 \quad FI = \sum_{i=1}^4 W_i \times VI_i \quad (2)$$

148 where VI_i is the importance score of the variable from selector i ($i=1,2,3,4$). W_i is the weighting factor
 149 based on prediction accuracy for each selection algorithm and determined by:

$$150 \quad W_i = \frac{MAPE_{\max} - MAPE_i}{MAPE_{\max} - MAPE_{\min}} \quad (3)$$

151 where $MAPE_i$ is the prediction error in terms of mean absolute percentage error based on the i^{th}
 152 selector. $MAPE_{\max}$ and $MAPE_{\min}$ are the maximum and minimum prediction errors from the four
 153 selectors, respectively. The definition of MAPE is shown in equation 4.

154

155 Backward elimination algorithm is applied to remove irrelevant and less important variables. The
 156 resulted variables are regarded as the input variable for training the next-level selectors. Fig. 5 shows
 157 the flow chart of the stacking framework in ensemble variable selection. The selection process repeats
 158 until a stop condition is met, either the prediction accuracy approaches the goal or the maximum
 159 number of epochs is reached. In this case, the optimal input variables for data-driven modelling of
 160 gas-liquid CO2 flow measurement are obtained.

161

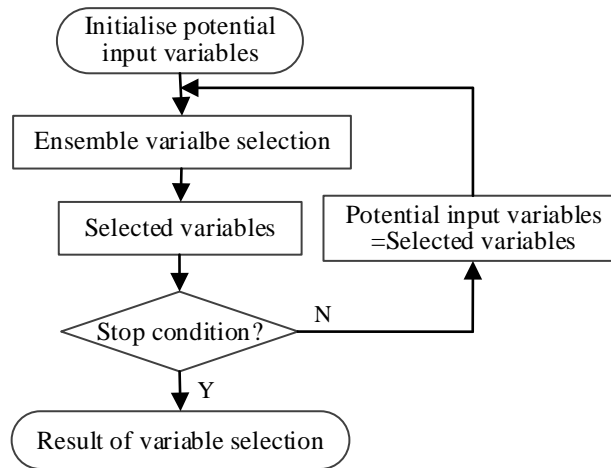


Fig.5 Block diagram of combination strategy based on stacking

162

163

164

165 3. Experimental results and discussion

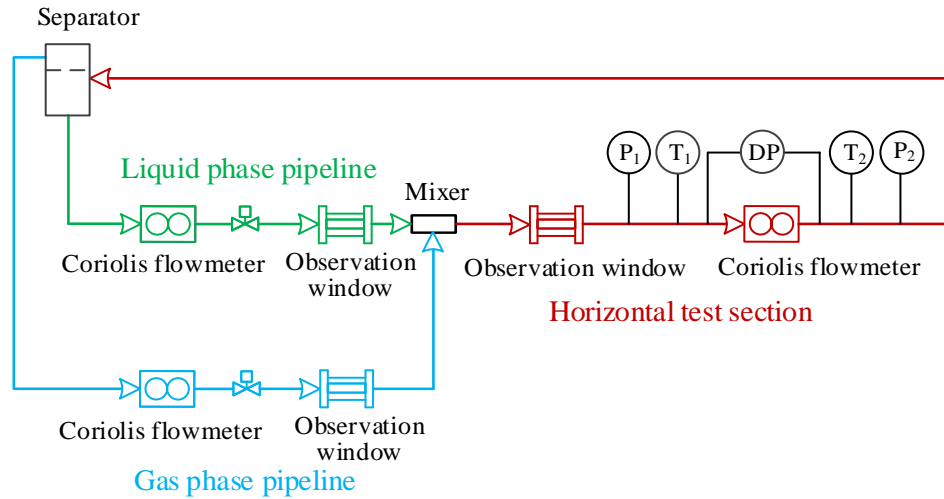
166 3.1 Experimental conditions

167 The test rig used in this study for gas-liquid two-phase CO₂ flow is shown in Fig. 6. The single gas
 168 phase and single liquid phase CO₂ flows are mixed through the mixer and form two-phase flow at the
 169 horizontal test section. During the mixing process, the temperature and pressure were kept around the
 170 gas/liquid transition line according to the phase diagram of CO₂ in order to achieve gas/liquid
 171 two-phase flow conditions. Meanwhile, the temperature and pressure were maintained at constant
 172 values via the control system to reduce the likelihood of state change during each test run. Meanwhile,
 173 the flowmeter under test was installed at only 1.6 m away from the gas-liquid mixer and the difference
 174 in temperature between the mixer and the test section was less than 1 °C, so the state change between
 175 gas CO₂ and liquid CO₂ was unlikely. Two separate Coriolis flowmeters were installed on the single
 176 liquid phase pipeline and single gas phase pipeline, respectively, to provide reference mass flowrate of
 177 CO₂ liquid and gas phases. The uncertainty of the flowmeters for the mass flowrate metering of liquid
 178 phase and gas phase are 0.16% and 0.35%, respectively. The accuracy of these meters is high enough
 179 to be used to obtain reliable reference values. At the test section, a Coriolis flowmeter (KROHNE
 180 OPTIMASS 6400 S15) is used as the target instrument to evaluate the proposed method. Pressure,
 181 temperature and differential-pressure transducers are added to capture additional information about the
 182 flow.

183

184 Experimental tests were conducted under the condition of total mass flowrate from 150~3500 kg/h

185 and the GVF from 1.11% ~ 88.44%. A total number of 541 experimental data were obtained. The time
 186 duration under each experimental condition is 100s. The temperature observed at the meter under test
 187 over all experiments ranged between 18°C and 25°C and the pressure was from 5.4 MPa ~ 6.5 MPa.
 188



189
 190 Fig. 6 Test rig for gas-liquid two-phase CO₂ flow
 191

192 All the variables derived directly or indirectly from sensor signals are listed in Table 1. $x1- x16$ are the
 193 original attributes collected by the sensors. As the data from the Coriolis flowmeter under test were
 194 updated around every 40 ms, $x1-x10$ were acquired via the General Device Concept (GDC) protocol
 195 with sampling rate of 48 Hz as per the sampling theorem. No filtering or limiting is applied to the data
 196 from the meter. Variables $x11-x16$ were acquired via an NI (National Instrument) data acquisition card
 197 from transducers with a sampling rate of 30 Hz. $x17- x30$ are the extended attributes including
 198 temperature difference, relative pressure difference and some statistical values of some original
 199 attributes. The extended attributes are regarded as potential variables, which contain useful
 200 information about the two-phase flow.

201
 202 Table 1 Input variables and their corresponding physical definitions

ID	Variable name	Physical definition
$x1$	Apparent mass flowrate (\dot{q})	The mass flowrate reading from the Coriolis flowmeter at the test section
$x2$	Process temperature (T)	The temperature reading from the Coriolis flowmeter at the test section
$x3$	Observed density (ρ_l)	The density reading from the Coriolis flowmeter at the test

		section
x4	Tube frequency (f)	The oscillation frequency reading from the Coriolis measuring tube inside the Coriolis flowmeter
x5	Two phase indicator	An indicator for the detection of a two-phase
x6	Time shift (t_d)	The time delay between the signals reading from the two motion sensors
x7	Observed flow velocity (v)	The flow velocity reading from the Coriolis flowmeter at the test section
x8	Sensor A level (V_A/V_{MAX})	The relative voltage amplitude of signals from the motion sensor A
x9	Sensor B level (V_B/V_{MAX})	The relative voltage amplitude of signals from the motion sensor B
x10	Drive level (I_D/I_{MAX})	The relative current amplitude of the driver output
x11	Inlet pressure (P_1)	The pressure of the fluid at the inlet of the Coriolis flowmeter
x12	Inlet temperature (T_1)	The temperature of the fluid at the inlet of the Coriolis flowmeter
x13	Outlet pressure (P_2)	The pressure of the fluid at the outlet of the Coriolis flowmeter
x14	Outlet temperature (T_2)	The temperature of the fluid at the outlet of the Coriolis flowmeter
x15	Temperature different(ΔT)	The temperature difference across the Coriolis flowmeter
x16	Differential pressure(DP)	The differential pressure across the Coriolis flowmeter
x17	Relative variance of DP	Variance/ <i>differential pressure</i>
x18	Sensor level different (ΔV)	The relative amplitude difference
x19	Pressure drop (DP/P_1)	Relative ratio of the pressure differential
x20	Damping ($x10/x8$)	Damping factor of the Coriolis measuring tubes
x21	Variance of flow velocity	The variance of the flow velocity
x22	Relative variance of flow velocity	Variance/ observed flow velocity
x23	Skewness of flow velocity	Skewness of flow velocity
x24	Variance of mass flowrate	Variance of mass flowrate
x25	Relative variance of mass flowrate	Variance/ apparent mass flowrate
x26	Skewness of mass flowrate	Skewness of mass flowrate
x27	Variance of density	Variance of density
x28	Relative variance of density	Variance/ <i>observed density</i>
x29	Skewness of density	Skewness of density
x30	Density drop ($(\rho_0-\rho_1)/\rho_0$)	Relative ratio of measuring section density to liquid density (ρ_0 is theoretical density of CO ₂ liquid phase at certain temperature and pressure)

203

204 3.2 Parameter optimization of individual selectors

205 To improve the performance of selectors there are several parameters need to be optimized. As for the

206 DT algorithm, the parameter tree depth is to be optimised. Bagging algorithm requires to determine
 207 the tree depth and the total number of training trees (i.e. model size). Apart from the tree depth and
 208 model size, GBDT and GBRF algorithms need to optimise the parameters of learning rate. As GBRF
 209 algorithm has the characteristics of random forest algorithm, the number of maximum input features is
 210 another parameter to be optimised. All the parameters are optimised through a trial-and-error
 211 approach.

212

213 During the process of parameter optimization, MAPE is used to evaluate the performance of each
 214 selector and defined as

$$215 \quad MAPE = \sum_{i=1}^n \left| \frac{\hat{Y}_i - Y}{Y} \right| \times \frac{100\%}{n} \quad (4)$$

216 where Y is the desired value of the target variable; \hat{Y}_i is the predicted value obtained by the i^{th} basic
 217 learner, and n is the number of all samples. All the MAPE values in the paper are the average
 218 prediction accuracy from 10-fold cross validation.

219

220 When the target variable is the mass flowrate of gas-liquid two-phase CO₂, the desired value Y is
 221 determined by q_m , the sum of liquid CO₂ mass flow (q_{ml}) and gas CO₂ mass flow (q_{mg}):

$$222 \quad q_m = q_{ml} + q_{mg} \quad (5)$$

223

224 When the target variable is the GVF of gas-liquid two-phase CO₂, the desired value Y is equal to α :

$$225 \quad \alpha = \frac{q_{vg}}{q_{vl} + q_{vg}} \times 100\% \quad (6)$$

226 where q_{vl} and q_{vg} are calculated volume flowrates of the liquid and gas phases, respectively.

227

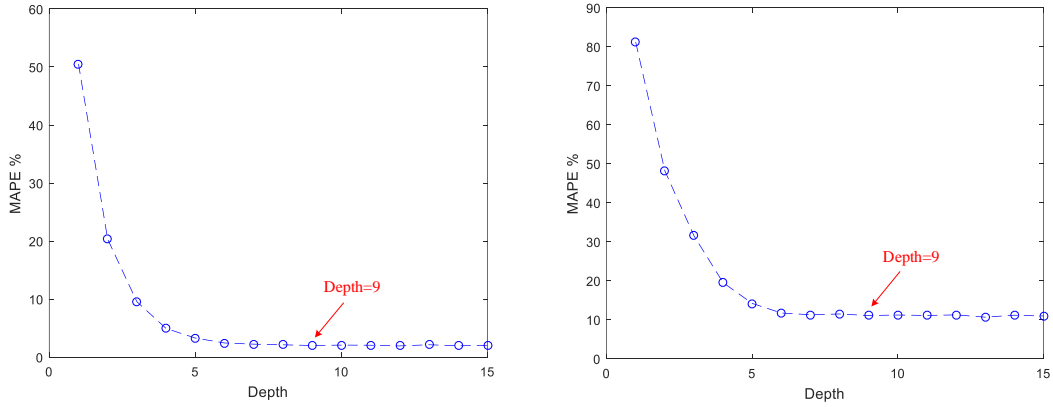
228 The reference value of GVF calculated from equation (6) is based on the assumption that there is no
 229 state change between the reference meters and the meter under test.

230

231 *1) Parameter optimisation for DT and Bagging models*

232 DT models are developed respectively with the tree depth from 1 to 15. As shown in Fig. 7(a) and 7(b),

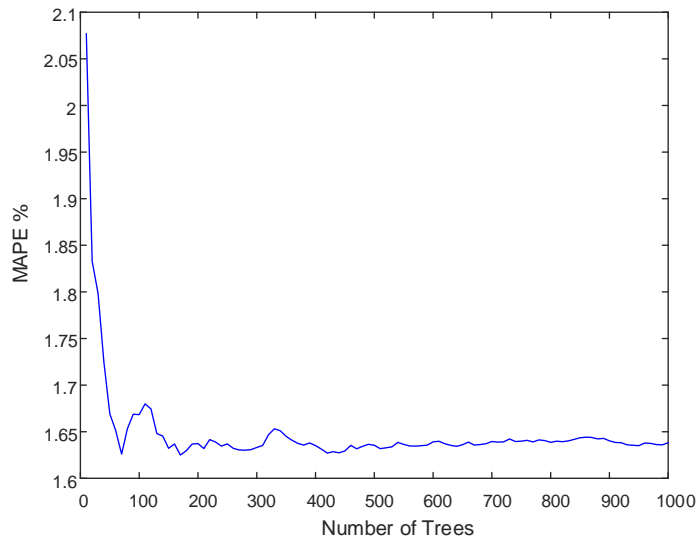
233 the MAPE decrease for both mass flowrate measurement and GVF prediction as the tree depth
 234 increases. When the depth is greater than 9, the prediction error does not change significantly. In order
 235 to avoid over-fitting, the model training depth is set to 9 in single decision tree models.



236
 237 (a) Performance of mass flowrate models (b) Performance of GVF models

238 Fig.7 Optimal depth in single tree models

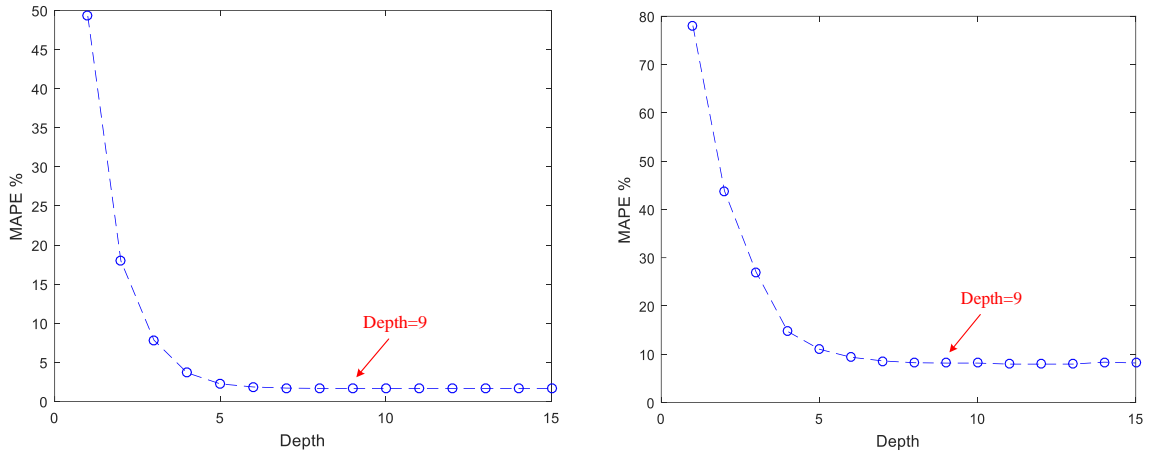
239
 240 As shown in Fig. 8, for the tree-based ensemble algorithms (Bagging, GBDT and GBRF), the
 241 prediction error of the target variable reaches a stable value (about 1.64% for mass flowrate) when the
 242 number of training trees is 500. Moreover, prediction error is no longer reduced as the number of
 243 training trees increases. Therefore, the number of training trees is set to 500 in the ensemble
 244 algorithms.



245
 246 Fig.8 Optimal number of trees in the tree-based ensemble algorithms

247
 248 Fig. 9 shows the process of determination of optimal depth in bagging algorithms which are used to

249 predict mass flowrate and GVF, respectively. In these bagging based models, the number of trees is set
 250 to 500. It is obvious that the MAPEs value of the prediction is approaching to the minimum of 1.64%
 251 for mass flowrate and 8.17% for GVF at a depth of 9. It also verifies the depth selection result in a
 252 single decision tree models. Compared with the prediction error of the single decision tree in Fig. 8,
 253 Bagging model with the same depth performs better as the result of ensemble learning.



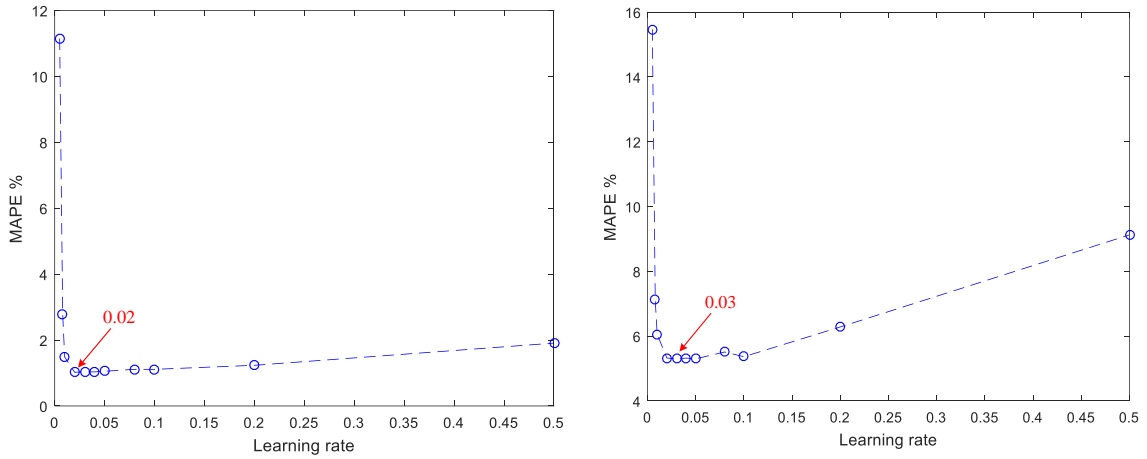
254 (a) Performance of mass flowrate models (b) Performance of GVF models

255 Fig.9 Optimal depth in Bagging algorithm

256
 257
 258 *2) Parameter optimisation for GBDT and GBRF models*

259 Different from bagging and RF, Gradient boosting (GB) is a serial ensemble and able to reduce bias
 260 and variance. GB often has low error values with stumps (a decision tree with a depth of 1) in deeper
 261 trees. Before determining their depth, it is necessary to determine the learning rate. At this time, the
 262 model is still trained with the depth of 9. The performance of GB models with respect to different
 263 learning rates is depicted in Fig.10. When the learning rate varies from 0.001~0.5, 500 trees are
 264 trained using GB algorithms at each learning rate. The prediction results of mass flowrate and GVF
 265 using GB method with different learning rates are shown in Fig. 11(a) and (b), respectively. When the
 266 learning rate is equal to 0.02, MAPE of mass flowrate models reaches the minimum value of 1.02%.
 267 For GVF models, the optimized learning rate is 0.03 to achieve a minimal MAPE of 5.3%.

268

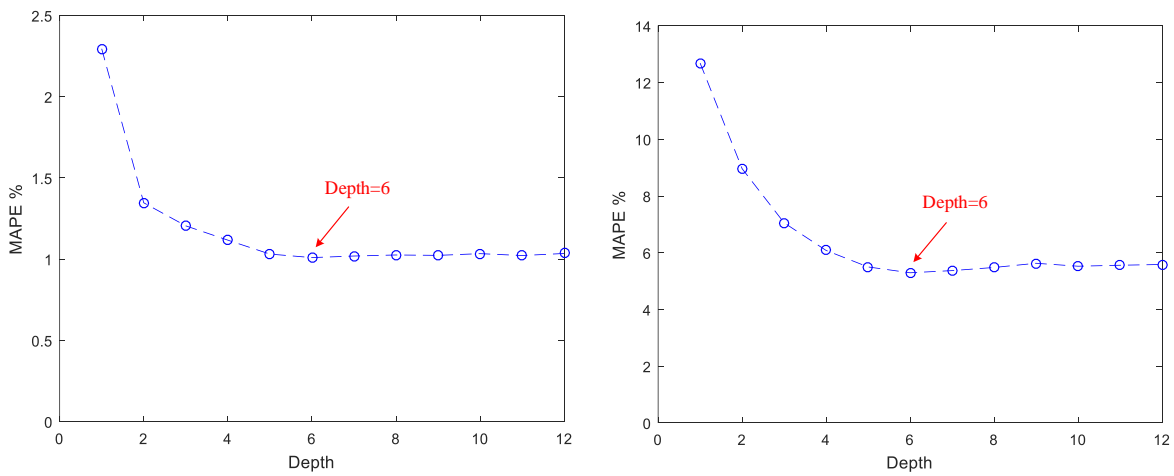


(a) Performance of mass flowrate models

(b) Performance of GVF models

Fig.10 Optimal learning rate in GB algorithm

Fig. 11 shows the performance of mass flowrate models and GVF models when the depth of the tree changes from 1 to 12 with a constant learning rate of 0.02 and 0.03. The MAPE value of models reaches the minimum of 1.03% for mass flowrate and 5.3% for GVF, when the depth of the tree is 6. For the serial gradient boosting based models, the model performance can be greatly improved by slightly increasing the depth of the decision tree due to the interaction among the potential variables. As there are some important variables with strong correlation in the potential input variables, it is necessary to apply the gradient lifting algorithm since both GBDT and GBRF algorithms play in a serial ensemble to the basic learners. They have the same depth parameter of 6.

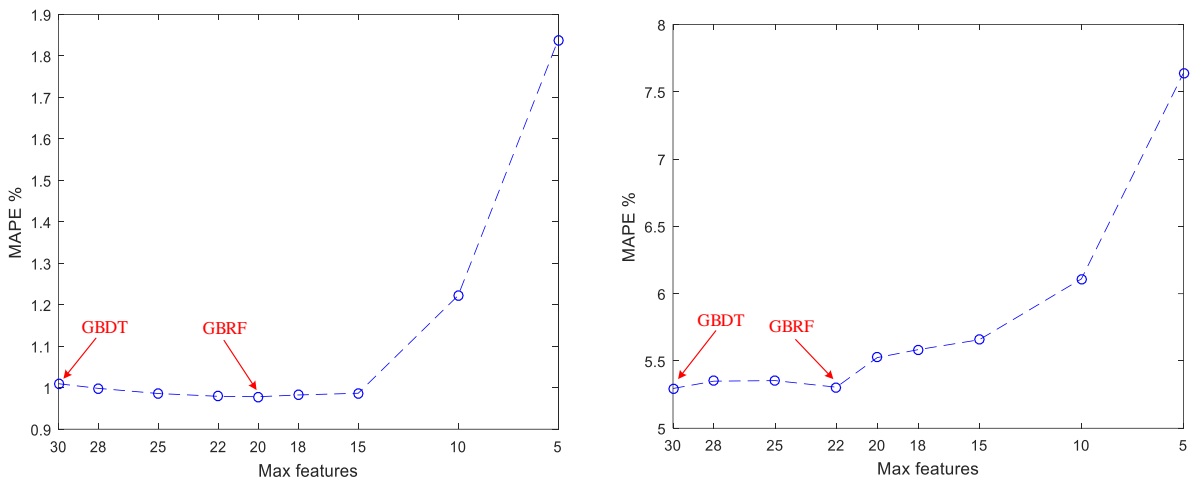


(a) Performance of mass flowrate models

(b) Performance of GVF models

Fig.11 Optimal depth in GB algorithm

285 Since GBRF is a combination of GB and random forest (RF) algorithms, it also need to select the
 286 maximum number of features (maxFeature) for model training. The maxFeature means the number of
 287 input variables. The learning rate and depth of GBRF are set to the determined values, respectively. As
 288 shown in Fig. 12, for the measurement of mass flowrate, the prediction error of the GBRF model
 289 reaches the minimum 0.98% when the maximum number of input features is 20. For the prediction of
 290 GVF, the GBRF model produce minimum MAPE of 5.3% when the maximum number of input
 291 features is 22. However, the optimal maxFeature for GBDT algorithm is 30 which is the total number
 292 of all possible input variables.



293 (a) Performance of mass flowrate models (b) Performance of GVF models
 294

295 Fig.12 Optimal max features in GBDT and GBRF algorithms

296 3) Summary of optimal parameters

297 According to the above analysis, the optimal parameters of the tree-based algorithms are obtained and
 298 summarized in Table 2.

299 Table 2 Optimal parameter of the basic learners

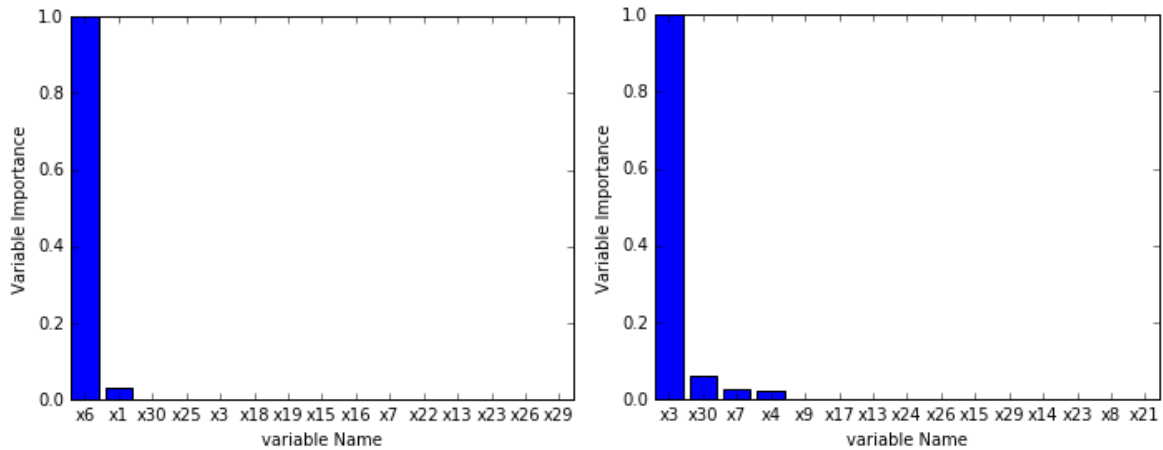
Parameter		DT	Bagging	GBDT	GBRF
Depth		9	9	6	6
The number of trees		1	500	500	500
Max number of input variables	Mass flowrate	30	30	30	20
	GVF	30	30	30	22
Learning rate	Mass flowrate	/	/	0.02	0.02
	GVF	/	/	0.03	0.03

300 *3.3 Implementation of ensemble variable selection*

301 Each of the algorithms is implemented based on the optimal parameters outlined in Table 2. The
 302 corresponding results of variable sorting and variable importance scores are obtained. The relative
 303 importance of the tree model is represented by the reduction in impurities due to the split on a specific
 304 variable set. For ensembles, the metric is averaged over the collection of base learners (Zhang et al.,
 305 2017), weighted average is used to derive the combined variable importance.

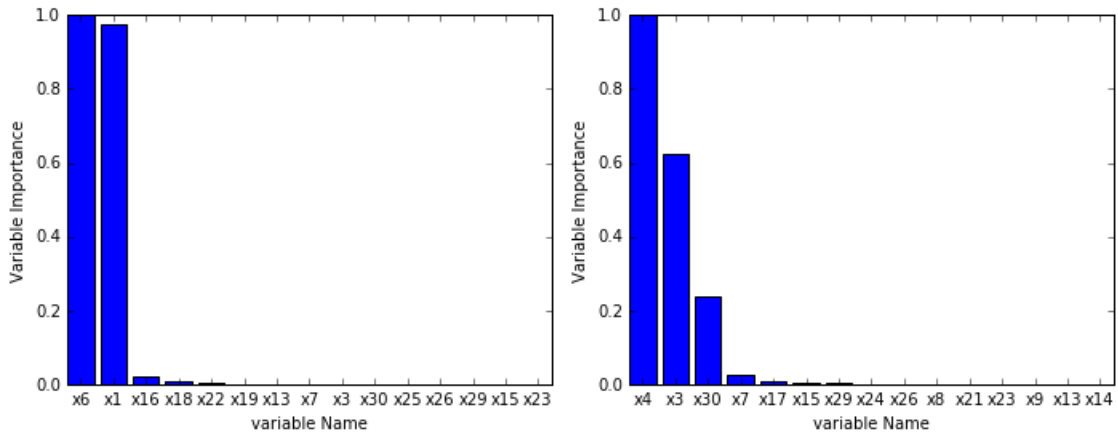
306

307 Fig.13 shows the order of first 15 variables obtained by each algorithm according to the normalized
 308 variable importance score. It can be seen from each figure, when the number of preselected input
 309 variables is 15, the variable importance of each algorithm varies. It is essential to effectively fuse and
 310 trim the sorting results of different algorithms to obtain the most accurate and concise set of input
 311 variables.



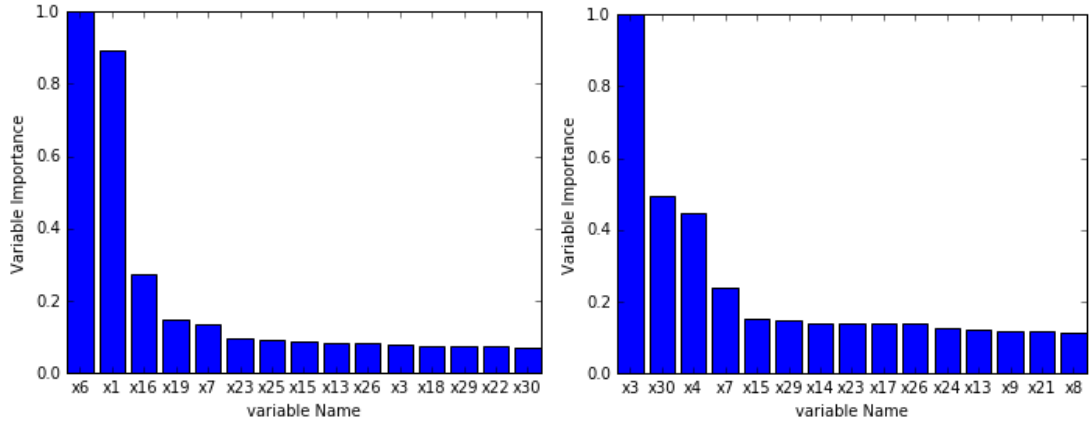
312

313 (a) Variable importance to mass flowrate using DT (b) Variable importance to GVF using DT



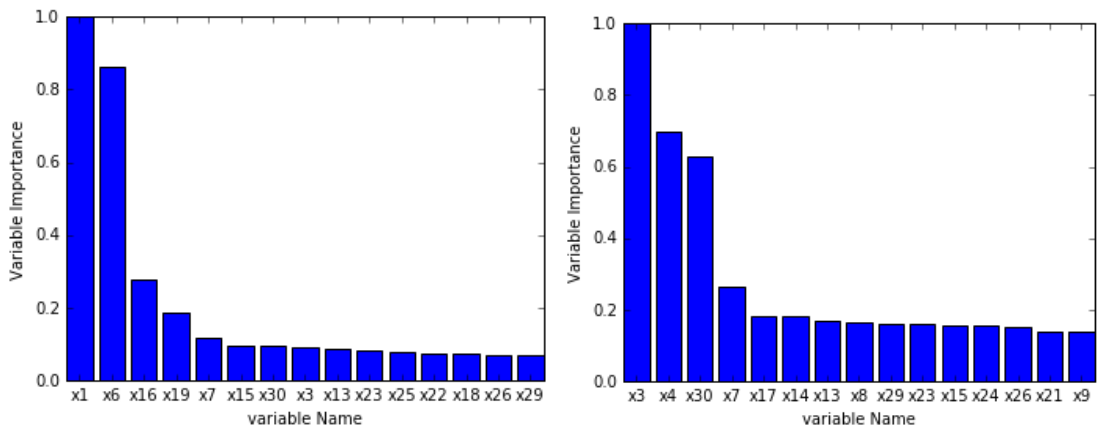
314

315 (c) Variable importance to mass flowrate using Bagging (d) Variable importance to GVF using
 316 Bagging



317

318 (e) Variable importance to mass flowrate using GBDT (f) Variable importance to GVF using GBDT



319

320 (g) Variable importance to mass flowrate using GBRF (h) Variable importance to GVF using GBRF

321 Fig.13 Variable importance to mass flowrate and GVF based on basics selectors

322

323 As shown in Fig. 13(a) and (b), the single algorithm DT can only select the first few most important
 324 variables, but cannot distinguish the importance of other variables. In Fig. 13 (c) and (d), the Bagging
 325 algorithm is an ensemble of 500 trees and the model performance is the average result of 500 trees.
 326 The algorithm can improve the polarization of importance and select more important variables.
 327 Therefore, prediction accuracy based on Bagging models is improved compared to single tree
 328 algorithm. In Fig. 13 (e) and (f), the use of GBDT algorithm is effective to get rid of the polarization
 329 phenomenon. The importance of all variables is smoothly changed. Therefore, the algorithm is
 330 superior to DT and Bagging algorithms, and has higher prediction accuracy. In Fig.13 (g) and (h),
 331 GBRF algorithm is a combination of gradient lifting and random forest algorithm, which further
 332 narrows the gap of the importance score of variables and further improves the prediction accuracy of
 333 target variables.

334
 335
 336
 337
 338
 339
 340
 341

Tables 3 and 4 summarize the results of the combined importance scores of the potential input variables. The algorithm performs 13 iterations of variable ranking and 12 iterations of variable selection. Backward elimination algorithm is applied to remove less important variables. For different prediction target, the results of variable selection are quite different. Finally, the validity of the selected variables is verified through assessing the performance of data-driven models.

Table 3 Results of variable selection at all levels for mass flowrate

The order of importance of tree-based feature selection													
Number Index	30	25	20	15	10	8	7	6	5	4	3	2	1
<i>x1</i>	<i>x1</i>	<i>x6</i>	<i>x6</i>	<i>x1</i>	<i>x6</i>	<i>x6</i>	<i>x6</i>	<i>x6</i>	<i>x6</i>	<i>x1</i>	<i>x1</i>	<i>x6</i>	<i>x6</i>
<i>x2</i>	<i>x6</i>	<i>x1</i>	<i>x1</i>	<i>x6</i>	<i>x1</i>	<i>x1</i>	<i>x1</i>	<i>x1</i>	<i>x1</i>	<i>x6</i>	<i>x6</i>	<i>x1</i>	
<i>x3</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>	<i>x16</i>		
<i>x4</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>	<i>x19</i>			
<i>x5</i>	<i>x7</i>	<i>x25</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>				
<i>x6</i>	<i>x25</i>	<i>x9</i>	<i>x25</i>	<i>x13</i>	<i>x25</i>	<i>x13</i>	<i>x25</i>	<i>x25</i>					
<i>x7</i>	<i>x22</i>	<i>x7</i>	<i>x15</i>	<i>x25</i>	<i>x13</i>	<i>x25</i>	<i>x13</i>						
<i>x8</i>	<i>x23</i>	<i>x22</i>	<i>x30</i>	<i>x15</i>	<i>x15</i>	<i>x15</i>							
<i>x9</i>	<i>x26</i>	<i>x23</i>	<i>x3</i>	<i>x23</i>	<i>x22</i>								
<i>x10</i>	<i>x29</i>	<i>x29</i>	<i>x26</i>	<i>x22</i>	<i>x23</i>								
<i>x11</i>	<i>x15</i>	<i>x30</i>	<i>x22</i>	<i>x30</i>									
<i>x12</i>	<i>x21</i>	<i>x24</i>	<i>x23</i>	<i>x26</i>									
<i>x13</i>	<i>x18</i>	<i>x26</i>	<i>x13</i>	<i>x18</i>									
<i>x14</i>	<i>x24</i>	<i>x15</i>	<i>x18</i>	<i>x29</i>									
<i>x15</i>	<i>x8</i>	<i>x18</i>	<i>x29</i>	<i>x3</i>									
<i>x16</i>	<i>x28</i>	<i>x8</i>	<i>x21</i>										
<i>x17</i>	<i>x27</i>	<i>x3</i>	<i>x2</i>										
<i>x18</i>	<i>x17</i>	<i>x13</i>	<i>x24</i>										

<i>x19</i>	<i>x9</i>	<i>x21</i>	<i>x9</i>										
<i>x20</i>	<i>x30</i>	<i>x2</i>	<i>x8</i>										
<i>x21</i>	<i>x2</i>	<i>x28</i>											
<i>x22</i>	<i>x3</i>	<i>x17</i>											
<i>x23</i>	<i>x13</i>	<i>x20</i>											
<i>x24</i>	<i>x12</i>	<i>x12</i>											
<i>x25</i>	<i>x20</i>	<i>x27</i>											
<i>x26</i>	<i>x4</i>												
<i>x27</i>	<i>x14</i>												
<i>x28</i>	<i>x11</i>												
<i>x29</i>	<i>x10</i>												
<i>x30</i>	<i>x5</i>												

342

343

Table 4 Results of variable selection at all levels for GVF

The order of importance of tree-based feature selection													
Number \ Index	30	25	20	15	10	8	7	6	5	4	3	2	1
<i>x1</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>	<i>x3</i>
<i>x2</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>	<i>x30</i>
<i>x3</i>	<i>x4</i>	<i>x4</i>	<i>x4</i>	<i>x4</i>	<i>x4</i>	<i>x4</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	
<i>x4</i>	<i>x29</i>	<i>x23</i>	<i>x29</i>	<i>x7</i>	<i>x7</i>	<i>x7</i>	<i>x4</i>	<i>x4</i>	<i>x14</i>	<i>x14</i>			
<i>x5</i>	<i>x23</i>	<i>x29</i>	<i>x23</i>	<i>x29</i>	<i>x14</i>	<i>x24</i>	<i>x14</i>	<i>x14</i>	<i>x4</i>				
<i>x6</i>	<i>x26</i>	<i>x26</i>	<i>x26</i>	<i>x24</i>	<i>x24</i>	<i>x14</i>	<i>x24</i>	<i>x24</i>					
<i>x7</i>	<i>x24</i>	<i>x15</i>	<i>x15</i>	<i>x26</i>	<i>x23</i>	<i>x23</i>	<i>x23</i>						
<i>x8</i>	<i>x15</i>	<i>x24</i>	<i>x24</i>	<i>x14</i>	<i>x29</i>	<i>x29</i>							
<i>x9</i>	<i>x17</i>	<i>x8</i>	<i>x17</i>	<i>x23</i>	<i>x26</i>								
<i>x10</i>	<i>x21</i>	<i>x17</i>	<i>x13</i>	<i>x15</i>	<i>x15</i>								
<i>x11</i>	<i>x8</i>	<i>x9</i>	<i>x7</i>	<i>x17</i>									
<i>x12</i>	<i>x9</i>	<i>x18</i>	<i>x8</i>	<i>x8</i>									

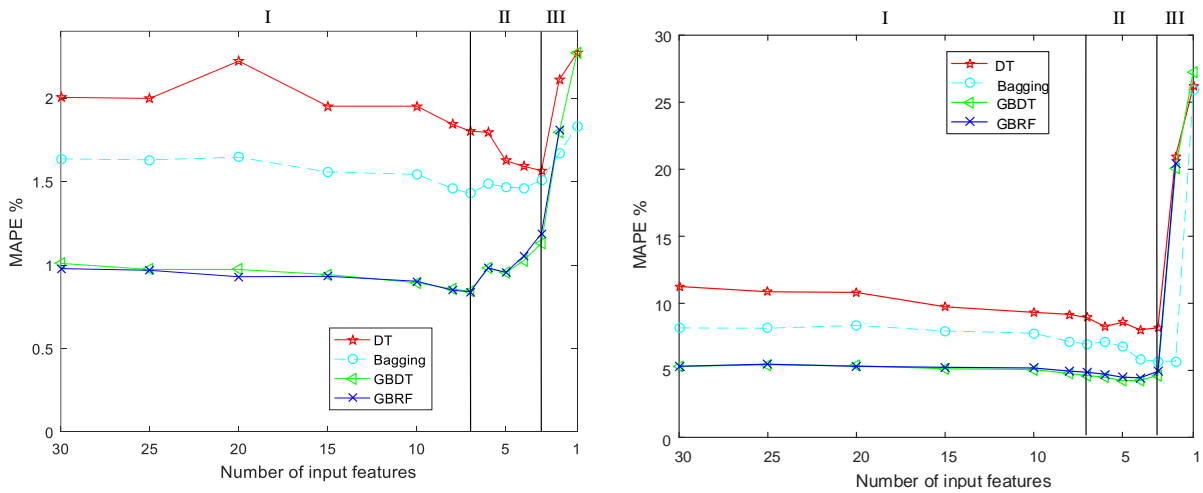
<i>x13</i>	<i>x27</i>	<i>x27</i>	<i>x14</i>	<i>x13</i>									
<i>x14</i>	<i>x18</i>	<i>x21</i>	<i>x21</i>	<i>x9</i>									
<i>x15</i>	<i>x7</i>	<i>x7</i>	<i>x9</i>	<i>x21</i>									
<i>x16</i>	<i>x28</i>	<i>x25</i>	<i>x27</i>										
<i>x17</i>	<i>x14</i>	<i>x14</i>	<i>x18</i>										
<i>x18</i>	<i>x25</i>	<i>x28</i>	<i>x25</i>										
<i>x19</i>	<i>x11</i>	<i>x13</i>	<i>x28</i>										
<i>x20</i>	<i>x19</i>	<i>x16</i>	<i>x16</i>										
<i>x21</i>	<i>x12</i>	<i>x22</i>											
<i>x22</i>	<i>x13</i>	<i>x12</i>											
<i>x23</i>	<i>x22</i>	<i>x11</i>											
<i>x24</i>	<i>x2</i>	<i>x19</i>											
<i>x25</i>	<i>x16</i>	<i>x2</i>											
<i>x26</i>	<i>x20</i>												
<i>x27</i>	<i>x5</i>												
<i>x28</i>	<i>x10</i>												
<i>x29</i>	<i>x6</i>												
<i>x30</i>	<i>x1</i>												

344

345 In the process of variable selection, the sooner the feature is removed, the less important it is. As
346 shown in Fig.10, this paper divides the elimination process of feature selection into three stages. At the
347 first stage, the range of input features retained after feature selection 7~30. At this stage, unimportant
348 variables are gradually eliminated until 7 variables are retained. Because the removed variables are
349 insignificant to the target variable, the prediction accuracy of the model is gradually improved. At the
350 second stage, the range of input variables retains after feature selection is 3~7. At this stage, it can be
351 seen that the lack of sub-important variables leads to a small increase in the prediction error, but it
352 does not have much impact on the overall performance. As shown in Fig.14(a), the prediction
353 performance of mass flowrate models is to be improved when the number of selected features is 5.
354 This is because the extended variables of important variables play an important role in predicting mass

355 flowrate. Fig. 14(b) shows the performance of GVF models. The MAPE value tends to increase when
 356 the number of selected features is more than 3. At the third stage, the range of input features retained
 357 after feature selection is 1~3. At this stage, one of the important variables required by the model is
 358 removed by fusion method, which leads to a rapid decline in the prediction accuracy of the model.
 359 Therefore, the optimal number of input variables for the prediction model should be determined at the
 360 second stage.

361



362

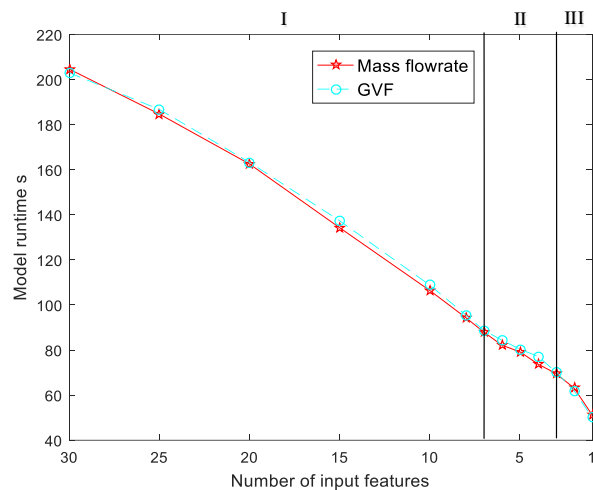
363 (a) CO₂ mass flowrate models

363 (b) CO₂ GVF models

364

364 Fig.14 Performance comparison of variable selection at different levels

365



366

367 Fig.15 Model efficiency comparison of mass flowrate and GVF at different levels

368

369 The model runtime in Fig. 15 is obtained, which represents the average run time of 54 test data (The
 370 processor of the computer used is Intel Xeon E5-2640 V4 CPU @ 2.4 GHz). The total running time of

371 each iteration is the sum of the running time of the four selection algorithms and their combination
 372 time. As shown in Fig. 16, as the number of input features is gradually reduced, the complexity of the
 373 model is reduced, and the time efficiency of the model is greatly improved. When the number of input
 374 features of the prediction model is 4, the efficiency of the GBRF algorithm is increased by 3 times
 375 compared with the feature number of 30. When the feature selection process reaches the second stage,
 376 the removal of a large number of variables unrelated to the predicted target value improves the model
 377 efficiency rapidly. This point indicates that the fusion variable selection method in this paper is very
 378 effective for large data measurement systems with a large number of input variables. Under the same
 379 number of input features, the running time of GVF prediction model is slightly longer than that of
 380 mass prediction model. This is because the relationship between the target variable GVF and the
 381 variable candidates is more complex and difficult to find.

382

383 The ensemble tree algorithm can not only select variables, but also derive the predicted results of
 384 target variables. Three indicators are integrated to determine the final variable selection results:
 385 prediction accuracy (MAPE), model efficiency (running time) and model complexity (number of input
 386 variables) in Figs 14 and 15. In Fig.14 (a), the MAPE of the three ensemble algorithms reaches the
 387 minimum value when seven variables are used as model inputs. In Fig. 14 (b), the MAPE of the three
 388 ensemble algorithms reaches the minimum value when four variables are used as model inputs. Based
 389 on the analysis results of the above three evaluation indicators, four input variables are determined as
 390 the final selection results for both mass flowrate and GVF prediction models. The results of variable
 391 selection are shown in table 5. The MAPE value increases slightly for mass flowrate when the number
 392 of model input variables is 4. However, the model efficiency is greatly improved and the model
 393 complexity is reduced.

394

395 The input variable for the mass flowrate prediction model is determined as $\{x1\ x6\ x16\ x19\}$ while the
 396 input variable for the GVF prediction model is $\{x3\ x30\ x7\ x14\}$. The input variables selection results
 397 using the proposed the tree-based heterogeneous ensemble approach is summarized in Table 5.

398

399 Table 5 Result of selection of four input variables

Index	Mass flowrate	GVF
-------	---------------	-----

1	<i>x1</i> -Apparent mass flowrate (\dot{m})	<i>x3</i> -Observed density (ρ_1)
2	<i>x6</i> -Time shift (t_d)	<i>x30</i> -Density drop ($(\rho_0-\rho_1)/\rho_0$)
3	<i>x16</i> -Differential pressure (DP)	<i>x7</i> -Observed flow velocity (v)
4	<i>x19</i> - Pressure drop (DP/P_1)	<i>x14</i> -Outlet temperature (T_2)

400

401 According to the physical meaning of the input variables and previous theoretical and experimental
402 study (Henry et al., 2006; Li et al., 2018), the selection results are further analyzed. Apparent mass
403 flowrate (\dot{m}) *x1* is measured in the horizontal test section, even though the CMF produces large errors
404 in measuring the mass flowrate of two-phase flow (the original measurement error). It is still related to
405 the desired CO₂ mass flowrate. Time shift (t_d) *x6* and apparent mass flowrate has a functional
406 relationship, so time shift is also correlated with the desired CO₂ mass flowrate. As the liquid CO₂
407 flowing through the meter with various gas CO₂ entrainment, the pressure difference *x16* across the
408 Coriolis flowmeter and relative pressure difference *x19* can characterise the mixed CO₂ flow to some
409 extent.

410

411 For variable selection of GVF, *x3* is the observed density of gas-liquid mixture phase, and density drop
412 *x30* is derived from the observed density and the liquid phase density, which can somehow reflect
413 GVF. Although *x7* observed flow velocity is not accurate two-phase flow velocity, it still can reflect
414 the variation of mixed flow in the pipe. As the physical properties of CO₂ are very sensitive to the
415 variations in fluid temperature and pressure, so when temperature increases, phase change from liquid
416 to gas may occur and hence increasing GVF. Therefore, *x14* temperature is also an important variable
417 for GVF prediction.

418

419 In Table 5, *x1* (apparent mass flowrate), *x6* (time shift), *x16* (differential pressure) and *x19* (pressure
420 drop) are selected for mass flowrate measurement. *x3* (observed density), *x30* (density drop), *x7*
421 (observed flow velocity) and *x14* (outlet temperature) are selected for GVF prediction. Variables *x1* &
422 *x6*, *x3* & *x30*, and *x16* & *x19* look like highly redundant pairs. However, the mathematical relationship
423 between each pair is complex particularly in the case of gas-liquid two-phase flow. For instance, *x1* is
424 derived from *x6*, but their exact relationship depends on the fluid temperature and material properties
425 of the sensing tube (Wang et al 2017). Since *x1* includes temperature compensation and material
426 property effect, both *x1* and *x6* are selected in this case. *x19* is the ratio of differential pressure (*x16*) to

427 the inlet pressure (x_{11}). x_{30} is the relative ratio of the observed density (x_3) to liquid density ρ_0 (ρ_0 is
428 the theoretical density of CO₂ liquid phase at certain temperature and pressure). Additional fluid
429 information is included in x_6 , x_{19} and x_{30} than those in x_1 , x_{16} and x_3 .

430

431 The selection processes produce different importance scores for these variables. In consideration of the
432 prediction accuracy and model complexity of the data-driven models, the combination of the variables
433 outlined in table 5 are taken as the ‘optimal’ inputs to the tree-based models. If both variables from
434 each pair are used as input features, the MAPE values of the tree-based models will be reduced by at
435 least 0.5% for mass flowrate and 7% for GVF. In this case, these pairs provide complementary
436 information for the data-driven models.

437

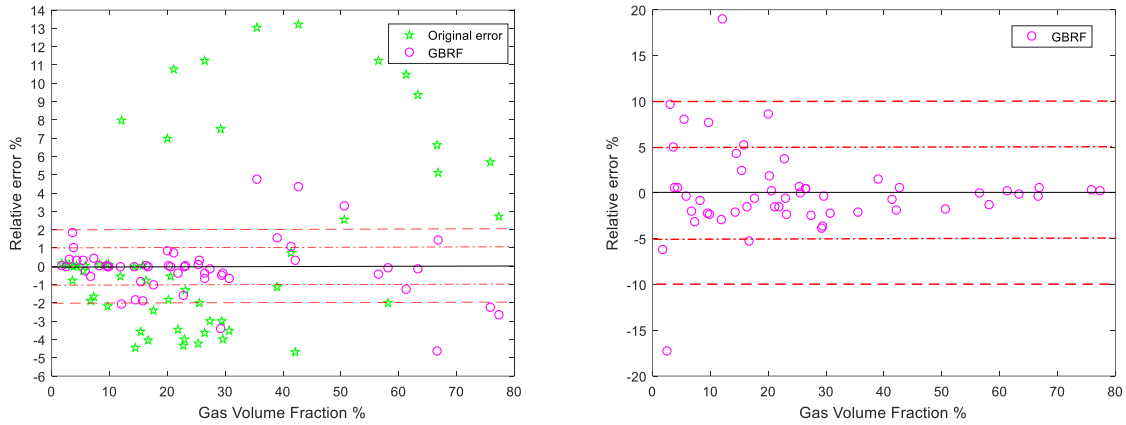
438 *3.4 Evaluation of ensemble variable selection*

439 GBRF models were developed based on the selected variables. The model performance was assessed
440 with 54 test samples (mass flowrate range: 212 kg/h~3449 kg/h and GVF range: 1.82%~77.29%).
441 As can be seen from Fig.16(a), the method proposed in this paper has improved the measurement
442 accuracy over the direct mass flow measurements. The prediction accuracy of GVF from the GBRF
443 model with selected input variables is shown in Fig.16(b). It turns out that the relative error in mass
444 flowrate measurement is mostly within $\pm 1.0\%$ and GVF prediction mostly within $\pm 5.0\%$.

445

446 When GVF is low, the gas entrainment has less effect on the vibration of the flowtubes in the Coriolis
447 mass flowmeter, which makes the apparent mass flowrate still very close to the true mass flowrate. In
448 this case, it is easier for the data-driven models to correct the errors. As GVF increases, the large
449 bubbles in the mixed flow lead to larger difference and nonlinear relationship between the apparent
450 mass flowrate and true mass flowrate. In this case, it is more challenging for the data-driven models to
451 derive the relationship and hence relative large errors than low GVF conditions.

452



(a) Mass flowrate

(b) GVF

Fig.16 Model performance with the selected variables

To further verify the selected results, a comparative experiment was conducted. A model is established with two input variables $x1$ (apparent mass flowrate) and $x30$ (density drop), respectively, for mass flow measurement and GVF prediction. A total of 108 samples for mass flowrate measurement and GVF prediction are tested. The performance of the models with two inputs and the selected inputs are summarised in Table 6. The results demonstrate that more test results from the model with selected inputs lie in the expected error range. 73% of test data produce relative error within $\pm 1.2\%$ for mass flow measurement and 80% of test data produce relative error within $\pm 4\%$ for GVF prediction.

Table 6 Performance comparison of data-driven models with different inputs

Comparison of mass flow prediction		
	$\{x1,x30\}$	$\{x1,x6,x16,x19\}$
Relative error within $\pm 1.2\%$	70%	73%
Comparison of GVF prediction		
	$\{x1,x30\}$	$\{x3,x7,x14,x30\}$
Relative error within $\pm 4\%$	69%	80%

4. Conclusions

A tree based heterogeneous embedded ensemble approach has been proposed for variable selection and applied to gas-liquid CO₂ two-phase flow measurement in this paper. Based on the combination

470 strategy of stacking and weighted averaging, the proposed method fuses the variable selection results
471 from four single selectors. At the same time, mass flowrate measurement and GVF prediction of
472 gas-liquid two-phase CO₂ flow have been carried out using the selected variables as inputs to the
473 GBRF models. The relative error of mass flowrate from the GBRF model is mostly within 1% with the
474 selected input variables (*apparent mass flow rate, time shift, differential pressure and pressure drop*).
475 The prediction error of GVF is mostly less than 5% using the selected input variables (*observed*
476 *density, density drop, observed flow velocity, outlet temperature*). The outcome from such modelling
477 research will help to enhance the understanding of two-phase flow measurement. Meanwhile, the
478 results presented in the paper demonstrate that the proposed heterogeneous ensemble approach is
479 capable of providing a small number of input variables and developing effective data-driven models
480 for multiphase flow measurement. In the further, more effort will be made to improve the
481 transferability of the developed data-driven model.

482

483 Engineering judgement here is still important as we have some knowledge of the two-phase flow and
484 Coriolis sensing process. Meanwhile, research is ongoing through analytical modelling of the
485 gas-liquid two-phase flow, which is a related area of research we are working on. The results from
486 such modelling research will help enhance engineering judgement. However, the variable selection as
487 reported in this paper will assist the optimisation of the machine learning models significantly. The
488 results presented in the paper demonstrate that the proposed heterogeneous ensemble approach is
489 capable of providing a small number of input variables and developing effective data-driven models
490 for multiphase flow measurement.

491

492 **Acknowledgements**

493 The authors would like to acknowledge the financial support of the National Natural Science
494 Foundation of China (No. 61973113 and No. 62073135). This work is also supported by Fundamental
495 Research Funds for the Central Universities (2020MS015) and by the UK CCS Research Centre
496 (www.ukccsrc.ac.uk). The UK CCSRC is funded by the EPSRC as part of the RCUK Energy
497 Programme.

498

499

500 **References**

- 501 [1] AL-Qutami, T. A., Ibrahim, R., Ismail, I., Ishak, M. A., 2018. Virtual multiphase flowmetering
502 using diverse neural network ensemble and adaptive simulated annealing. *Expert Syst. Appl.* 93,
503 72-85.
- 504 [2] Breiman, L., 1996. Stacked regressions. *Mach. Learn.* 24(1),49-64.
- 505 Efron, B., Hastie, T., Johnstone, I., Tibshirani, R., 2004. Least angle regression. *Ann. Stat.* 32(2),
506 407-499.
- 507 [3] Green, T., Reese, M., Henry, M., 2008. Two-phase CO₂ measurement and control in the Yates oil
508 field. *Measurement and Control* , 41 (7), 205-207.
- 509 [4] Henry, M., Tombs, M., Duta, M., Zhou, F.B., Mercado, R., Kenyery, F., Langansan, R., 2006.
510 Two-phase flow metering of heavy oil using a Coriolis mass flow meter: A case study. *Flow Meas.*
511 *Instrum.* 17 (6), 399-413.
- 512 [5] Kunze, J.W., Storm, R., Wang, T., 2014. Coriolis mass flow measurement with entrained gas.
513 *Sensors and Measuring Systems.* 17. ITG/GMA Symposium. VDE, 1-6.
- 514 [6] Li, J.L., Zhang, C.X., 2017. Ensembling variable selectors by stability selection for the Cox model.
515 In 2017 International Conference on Machine Learning and Cybernetics (ICMLC). 1, 35-41. IEEE.
- 516 [7] Li, M., Henry, M., Zhou, F., Tombs, M., 2019. Two-phase flow experiments with Coriolis Mass
517 Flow Metering using complex signal processing. *Flow Meas. Instrum.* 69 , 101613.
- 518 [8] Li, M., Henry, M., 2018. Complex signal processing for Coriolis mass flow metering in two-phase
519 flow, *Flow Meas. Instrum.* 64, 104-115.
- 520 [9] Mendes-Moreira, J., Soares, C., Jorge, A. M., Sousa, J. F. D., 2012. Ensemble approaches for
521 regression: A survey. *ACM Comput. Surv.* 45(1), 10.
- 522 [10] Nan, Y., Yang, Y., 2014. Variable selection diagnostics measures for high-dimensional
523 regression. *J. Comput. Graph. Stat.* 23(3),636-656.
- 524 [11] Shao, D., Yan, Y., Zhang, W.B, Sun, S.J, Sun, C.Y, Xu, L.J, 2020. Dynamic measurement of gas
525 volume fraction in a CO₂ pipeline through capacitive sensing and data driven modelling. *Int. J. of*
526 *Greenh. Gas Control.* 94, 102950.
- 527 [12] Tuv, E., Borisov, A., Runger, G., Torkkola, K. , 2009. Feature selection with ensembles, artificial
528 variables, and redundancy elimination. *J. Mach. Learn. Res.*10(Jul), 1341-1366.
- 529 [13] Wang, L.J, Yan, Y., Wang, X., Wang, T., 2017. Input variable selection for data-driven models of

- 530 Coriolis flowmeters for two-phase flow measurement. *Meas. Sci. Technol.* 28(3), 035305.
- 531 [14] Wang, L. J, Yan, Y., Wang, X., Wang, T., Duan, Q., Zhang, W., 2018. Mass flow measurement
532 of gas-liquid two-phase CO₂ in CCS transportation pipelines using Coriolis flowmeters. *Int. J. of*
533 *Greenh. Gas Control.* 68, 269-275.
- 534 [15] Wang, J.Z., Wu, L.S., Kong, J., Li, Y.X., Zhang, B.X., 2013. Maximum weight and minimum
535 redundancy: a novel framework for feature subset selection. *Pattern Recognit.* 46(6), 1616-1627.
- 536 [16] Xin, L., Zhu, M., 2012. Stochastic stepwise ensembles for variable selection. *J. Comput. Graph.*
537 *Stat.* 21(2), 275-294.
- 538 [17] Yan, Y., Wang, L. J, Wang, T., Wang, X., Hu, Y., Duan, Q., 2018. Application of softcomputing
539 techniques to multiphase flow measurement: A review. *Flow Meas. Instrum.* 60, 30-43.
- 540 [18] Zhang, C.X., Zhang, J.S., Yin, Q.Y., 2018. Early stopping aggregation in selectivevariable
541 selection ensembles for high-dimensional linear regression models. *Knowledge-Based Syst.*
542 153,1-11.
- 543 [19] Zhang, C.X., Wang, G.W., Liu, J.M., 2015. RandGA: injecting randomness into parallelgenetic
544 algorithm for variable selection. *J. Appl. Stat.* 42(3), 630-647.
- 545 [20] Zhang, C.X., Zhang, J.S., Yin, Q.Y., 2017. A ranking-based strategy to prune variable selection
546 ensembles. *Knowledge-Based Syst.* 125, 13-25.
- 547 [21] Zhu, M., Chipman, H.A., 2006. Darwinian evolution in parallel universes: A parallelgenetic
548 algorithm for variable selection. *Technometrics* 48(4), 491-502.
- 549 [22] Zhu, M., Fan, G., 2011. Variable selection by ensembles for the Cox model. *J. Stat. Comput.*
550 *Simul.* 81(12), 1983-1992.
- 551 [23] Zhou, Z.H., 2012. Ensemble methods: foundations and algorithms. Chapman and Hall/CRC.
- 552 [24] Zhang, W.B, Shao, D., Yan, Y., Liu, S., Wang, T., 2018. Experimental investigationsintothe
553 transient behaviours of CO₂ in a horizontal pipeline during flexible CCS operations. *Int. J. of*
554 *Greenh. Gas Control.* 79, 193-199.

555