# RSC Advances

1     # Modeling and Optimizing Performance of PVC/PVB

2     ## Ultrafiltration Membranes Using Supervised Learning

3     ## Approaches

4

5     Lina Chi [a,b, c*], Jie Wang [b], Tianshu Chu [b], Yingjia Qian [a], Zhenjiang Yu [a], Deyi Wu [a],

6     Zhenjia Zhang [a], Zheng Jiang [c], James O. Leckie [b]

7

8     *[a] School of Environmental Science and Engineering, Shanghai Jiao Tong University,*
9     *Shanghai, 200240, PRC*
10    *[b] The Center for Sustainable Development and Global Competitiveness (CSDGC),*
11    *Stanford University, Stanford, CA 94305, USA*
12    *[c] Faculty of Engineering and the Environment, University of Southampton,*
13    *Southampton, SO17 1BJ, UK*
14

15    * Correspondent author. Tel: +86 13816632156; E-mail address: <u>lnchi@sjtu.edu.cn</u>

16

17    **Abstract**

18    Mathematical models plays an important role in performance prediction and

19    optimization of ultrafiltration (UF) membranes fabricated via dry/wet phase inversion

20    in an efficient and economical manner. In this study, a systematic approach, namely, a

21    supervised, learning-based experimental data analytics framework, is developed to

22    model and optimize the flux and rejection rate of Poly (vinyl chloride) (PVC) and

23    Polyvinyl butyral (PVB) blend UF membranes. Four supervised learning (SL)

24    approaches, namely, the multiple additive regression tree (MART), the neural

25    network (NN), the linear regression (LR), and the support vector machine (SVM), are

26    employed in a rigorous fashion. The dependent variables representing membrane

27    performance response with regard to independent variables representing fabrication

28    conditions are systematically analyzed. By comparing the predicting indicators of the

29    four SL methods, the NN model is found to be superior to the other SL models with

30    training and testing R-squared values as high as 0.8897 and 0.6344, respectively, for

31    the rejection rate, and 0.9175 and 0.8093, respectively, for the flux. The optimal

32    combination of processing parameters and the most favorable flux and rejection rate

33    for PVC/PVB ultrafiltration membranes are further predicted by NN model and

34    verified by experiments. We hope the approach is able to shed light on how to

35    systematically analyzing multi-objective optimization issues for fabrication conditions

1

36   to obtain the desired ultrafiltration membrane performance based on complex

37   experiment data characteristics.

38

39   **Key words:** Poly (vinyl chloride) (PVC); Polyvinyl butyral (PVB); Supervised

40   Learning （SL）; Neural network (NN); modeling; Membrane fabrication

41   optimization

42

43   **1. Introduction**

44      Poly (vinyl chloride), or PVC, is commonly used to produce relatively

45   inexpensive ultrafiltration (UF) membranes due to its relative low cost, robust

46   mechanical strength, and other favorable physical and chemical properties, such as

47   abrasive resistance, acid and alkali resistance, microbial corrosion resistance, and

48   chemical performance stability[1]. Moreover, PVC membranes can usually maintain a

49   longer membrane life and remain intact after repeated cleaning with a wide variety of

50   chemical agents. However, the hydrophobic nature of PVC always leads to severe

51   fouling, thereby impeding its applications[1,2]. Thus, a critical challenge is to improve

52   the hydrophilicity of PVC membranes without interfering with their positive

53   characteristics so that PVC-based membranes can comply with industry requirements

54   for a wider range of applications.

55      In recent years, considerable research has been conducted in order to overcome

56   this problem. Among all available methods, polymer blends often exhibit superior

57   properties when compared with a standalone, individual component polymer; in

58   addition, the polymer blend method also has the advantages of a simple procedure for

59   preparation and easy control of physical properties for various compositional changes.

60   There are several polymers that have been studied as functional polymer pairs of PVC,

61   such as PMMA[1], PU[3], EVA[4], PEO[5], and PVB[6] among others. In most previous

62   studies[7,8], PVB is found to be one of the ideal polymers to blend with PVC due to its

63   well-predicted miscible properties, chemical similarity, and less unfavorable heat

64   while mixing. In addition, owing to the –OH bond, the PVC/PVB blend demonstrates

65   more hydrophilicity than the original PVC membrane [6,9] .

66      The selection of membrane material is essential for developing high-performance

67   membranes. However, due to the complexities of the fabrication process, even more

68   critical—especially when the membranes are made via a complex dry/wet phase

69 inversion—is a consistent and robust data analysis procedure for effectively analyzing

70 these membranes for better performance. Pure water flux (PWF) and rejection rate of

71 Bull Serum Albumin (BSA) are the most important performances for UF

72 membranes[10,11], depending not only upon the composition of the casting solution but

73 also upon the technical conditions used in the fabrication process. Typical variables of

74 importance for membrane development include the types and amounts of polymer,

75 additive, and the pore-forming agents used in the casting solution, the kind and

76 concentration of gelation medium, the evaporation time and temperature of the

77 spread-casting solution, the length of gelation period, and the temperature of gelation

78 bath[12] etc.. Some of the above mentioned variables have to be classified as categorical

79 variables, such as the type of the polymer, the pore-forming reagent, or the gelation

80 medium used, since they cannot be quantified. Remaining variables are quantitative

81 ones, including the temperature of evaporation or gelation, the amount of pore-

82 forming reagent added, and the duration of evaporation or gelation. Generally, these

83 complex influential factors in the membrane fabrication process would greatly delay

84 the development cycle and increase research and development (R&D) costs.

85 Therefore, it is worthwhile to investigate efficient statistical and computational

86 methods to optimize experiment design and to minimize the number of experiments.

87 Traditionally, statistically-based design of experiments (DOE) has been widely

88 used as a proper approach to optimize membrane parameters in membrane fabrication

89 processing[13-15]. However, DOE is based on the assumption that interactions between

90 factors are not likely to be significant[16,17], which is usually not the case in the real

91 world. When reducing the number of runs, a fractional factorial DOE becomes

92 insufficient to evaluate the impact of some of the factors independently[16]. Moreover,

93 it is also beyond the ability of DOE in dealing with categorical factors in experiments.

94 As a result, DOE has limitations in modeling a membrane fabrication process and in

95 optimizing the filtration performance of the membrane.

96 Recently, the supervised learning (SL) approach—a powerful method in

97 analyzing complex, but data-rich problems—has found strong application in diverse

98 engineering fields such as control, robotics, pattern recognition, forecasting, power

99 systems, manufacturing, optimization, and signal processing, etc. [18-20]. Although the

100 idea of solving engineering problems using SL has been around for decades, it has

101 been introduced only recently into the field of material studies[21]. There are several

102 publications discussing the application of SL to the modeling and optimization of

103  membrane fabrication. S. S. Madaeni modeled and optimized PES- and PS-membrane
104  fabrication using artificial neural networks[22], while Xi and Wang [23]reported that the
105  Support Vector Machine (SVM) model could be an efficient approach for optimizing
106  fabrication conditions of homemade VC-co-VAc-OH microfiltration membranes. Yet,
107  there are still a couple of key issues that need to be investigated. A systematic
108  framework for using SL approaches is required to discover the relationships between
109  membrane performance and complicated fabrication conditions.

110      The purpose of this research is to develop such a framework. More specifically,
111  we need first to evaluate experimental data quality, which is important in making
112  valid assumptions and selecting proper models for analyzing complex data. Secondly,
113  we need to develop an approach for efficiently employing reliable analysis models,
114  including the decision tree approach, neural network method, linear regression, and
115  support vector machine, for thoroughly analyzing all features and all responses of the
116  membranes, as opposed to current approaches that analyze only a single response with
117  regard to either one feature or all of the features. Finally, we need to select the most
118  suitable SL approach to predict the optimal combination of features for membrane
119  fabrication.

120

121  **2.  Experimental**

122  **2.1. Chemicals and materials**

123      Unless otherwise specified, all reagents and chemicals used were of analytical
124  grade. More specifically, PVC resin (Mw = $1.265 \times 10^5$ g/mol, and [η] = 240 mPa·s)
125  was supplied by Shanghai Chlor-Alkali Chemical Co., Ltd. Mw = $1.265 \times 10^5$ g/mol,
126  and [η] = 240 mPa·s. PVB (Mw = $3.026 \times 10^4$ g/mol, and [η] = 40 mPa·s) was
127  purchased  from  Tianjin  Bingfeng  Organic  Chemical  Co.,  Ltd..  N,N-
128  dimethylacetamide (DMAc) was purchased from Shanghai Lingfeng Chemical
129  Reagent Co., Ltd. PEG 600, PVP K90, and Ca(NO3)2 were purchased from Aladdin
130  Industrial Inc. BSA(Mw=67,000 g/mol) was supplied by Shanghai Huamei Biological
131  Engineering Company.

132  **2.2 Membrane fabrication**

133      PVC/PVB composite membranes were prepared by the non-solvent induced
134  phase inversion. The casting solutions, containing PVC, PVB, DMAc, and additives,
135  were prepared in a 250 mL conical flask and heated to approximately 30-80 °C in a
136  water bath while being stirred at 600 rpm using a digital stirring machine (Fluko, GE).

137   After the polymers had been dissolved completely and stirred for at least 24 h, the

138   resulting solution was degassed for at least 30 min until no gas bubbles were visible.

139   The solution was cast on a glass plate using an 8-inch wide doctor blade with a gap of

140   200 μm between the glass plate and blade. The temperature of the blade and the glass

141   plate was controlled between 30-80 °C. After a predetermined evaporation period,

142   ranging from 5 to 120 seconds, the film was immersed in a pure water or DMAc (with

143   volume concentration ranging from 10–80%) gelation bath maintained at 20°C. The

144   film was then removed from the glass plate and leached overnight in water in order to

145   completely remove any traces of solvent. Table S1 listed the various combination of

146   composition of casting solutions and corresponding processing parameters.

147   **2.3 Membrane characterization**

148       The pure water flux of the PVC/PVB blend ultrafiltration membranes was

149   measured at a temperature of 25 °C and under an operating pressure of 0.1 MPa after

150   pre-operating for 30 min. The flux of permeate was calculated according to Eq.(1):

151   $$J_w = V/(A \cdot t) \tag{1}$$

152   where $J_w$ (L/(m$^2$ • hr)) is the pure water flux, V (L) is the volume of the collected

153   permeate, and A (m$^2$)is the area of the membrane. In our study, the effective

154   membrane area is 0.0342 m$^2$ and t (hr) is the separation time.

155       Membrane retention ability was tested using 100 mg/L BSA at a temperature of

156   20 °C and under an operating pressure of 0.1 MPa. The concentrations of both the

157   feed water and the permeation water were determined using an ultraviolet

158   spectrophotometer (TU-1810, Beijing Purkinje Genera, China) at a wavelength of 280

159   nm. The percentage of the observed rejection solutes BSA phosphate buffer for each

160   permeate collected was calculated as the following Eq.(2):

161   $$R = (1 - C_p/C_f) \times 100\% \tag{2}$$

162   where $C_p$ is the permeate concentration and $C_f$ is the feed concentration.

163   **3. Analyzing membrane performance by SL approaches**

164       In both this section and in Section 4, we describe a systematic framework for

165   modeling and optimizing performance of PVC/PVB ultrafiltration membranes using

166   supervised learning approaches, consisting of the following: (1) methods for

167   analyzing raw datasets and their dependencies, (2) a general procedure and algorithms

168   of SL-based data processing, (3) detailed results analysis and comparisons among all

169    SL approaches, and (4) selection of the best learning approach for optimally

170    predicting experimental performance for analyzing membrane performance.

171    **3.1 Data structures and characteristics**

172        To better understand the potential inherent structures among independent and

173    dependent variables, in this section, we first describe data structures and

174    characteristics of experimental data sets. As listed in Table S1, there are a total of 68

175    valid experimental measurements. For each measurement, we have initially identified

176    and employed 9 processing parameters that are regarded as independent variables and

177    2 performance indicators that are regarded as dependent variables. Specifically, the

178    processing parameters are PVC Wt%, DMAc Wt%, Additive Wt%, Additive type

179    (PEG600, PVPk90, $Ca(NO_3)_2$), Casting solution temperature (℃), Evaporation time

180    (sec), Blade temperature (℃), Gelation bath type (Water, DMAc), and Bath

181    concentration (solute concentration in gelation bath) (mg/L). Note that the types of

182    additives and the gelation bath are categorical variables. The performance indicators,

183    including the rejection rate of BSA (%) and the flux ($L/(m^2·h)$), are numerical

184    variables. Through our preconditioning analysis, we find that the Wt% of polymers

185    and the Wt% of PVB have to be removed from the processing parameters because

186    they are dependent on, and correlated with, the change of PVC Wt%, DMAc Wt%,

187    and Additive Wt%. We introduce $k$ as the ratio of PVC Wt%/ Polymer Wt%, giving

188    us $0 < k < 1$. There exist following relationships:

189                PVC Wt%/k=Polymer Wt%                                    (3)

190                DMAc wt%+Polymer Wt%+Additive Wt%=100%         (4)

191        Before the data analysis process, we briefly verify the characteristics of the data

192    by scattering the measurement points under different parameter-indicator pairs in Fig.

193    1. If the processing parameters are categorical, box-plots are used instead of scatter

194    plots. Obviously, the rejection rate and the flux are negatively correlated. For

195    numerical parameters, PVC Wt% and DMAc Wt% have the strongest correlations

196    with flux and rejection rate, respectively, while evaporation time and blade

197    temperature have cross-like scatterings, thus indicating very weak correlations. Both

198    categorical parameters can provide considerable information for performance

199    prediction. This is especially true for the additive type, where the significant

200    differences of indicators are shown between different groups of additives. In general,

201    useful information can be found in the data for performance prediction, but there are

202    not enough measurements to estimate how the indicators are distributed with regard to

203    processing parameters. In other words, our predicted indicators using SL tools will

204    have a low bias but high variance, and we need to carefully balance the accuracy and
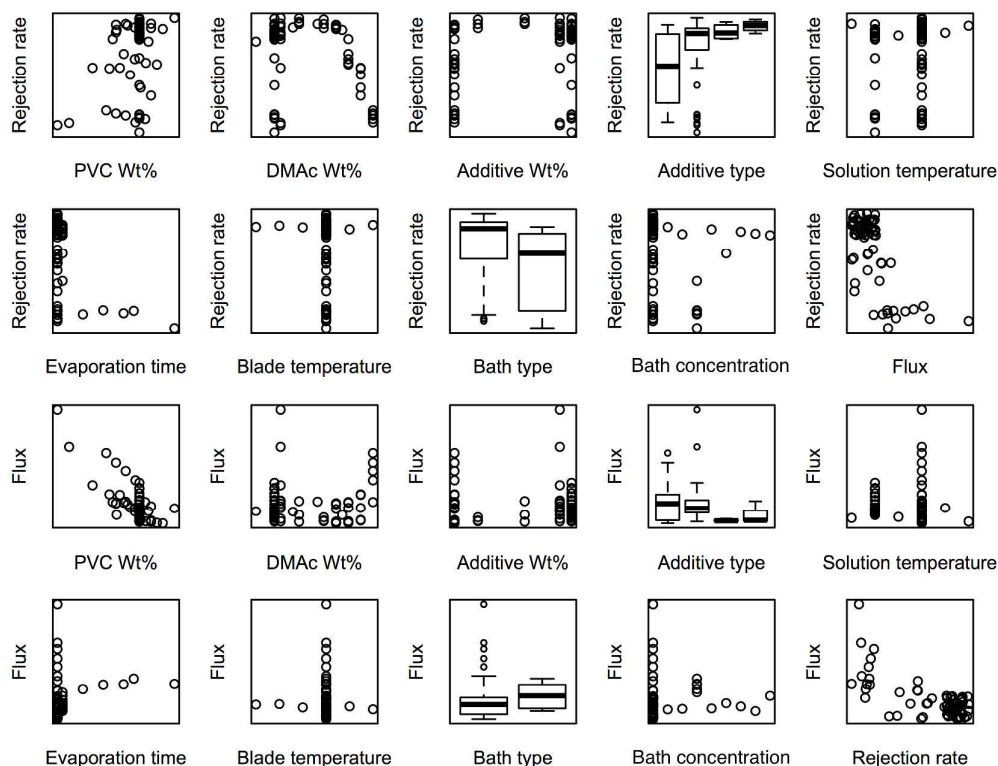
205    stability of modeling.



206
207                    Fig. 1 Scatter plots over measurements.
208

### 3.2 Supervised learning and data analysis procedures

### 3.2.1 General description and criteria

211    Different SL algorithms, including linear regression (LR), a multiple additive

212    regression tree (MART), a neural network (NN), and a support vector machine

213    (SVM), were introduced and implemented to find the potential influence of

214    processing parameters (predictors) on performance indicators (responses). The

215    advantages, limitations and assumptions when utilizing each SL algorithm were

216    described in Supporting information. To analyze the results, we train each SL

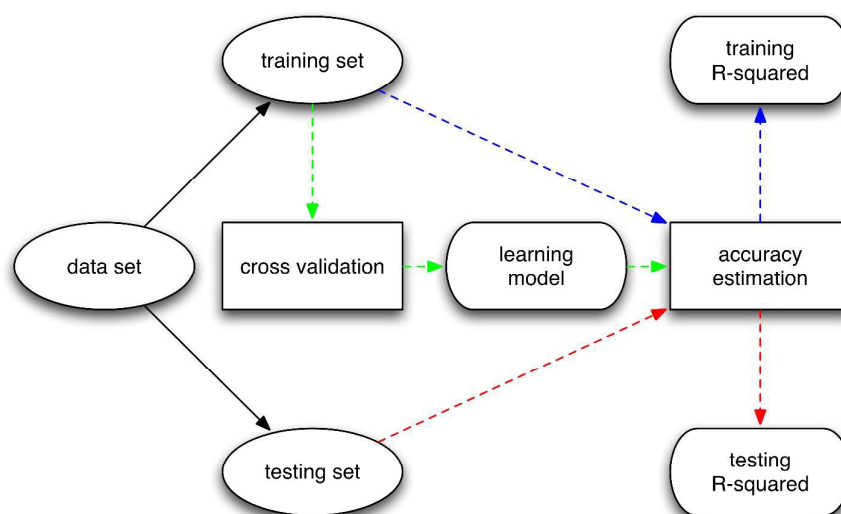217    algorithm over the whole data.

218    Furthermore, to estimate the accuracy of each SL algorithm, we apply the Monte

219    Carlo method by repeating the learning processes 50 times on our measurement data.

220    During each learning process, we first randomly split the data into a training set and a

221    testing set, with the ratio 50/18. Next, we train each SL model based on the predictors

222    of the training set with cross-validation and make predictions of responses over the

7

223    training and testing sets using the trained learning model. Finally, we estimate the

224    accuracy of each model by R-squared over the training and testing sets, computed as:

225    $$R^2 = 1 - \frac{\sum_{i=1}^{m}(\hat{y}^{(i)} - y^{(i)})^2}{\sum_{i=1}^{m}(\bar{y} - y^{(i)})^2} \qquad (5)$$

226    where m denotes the size of the data over which we perform predictions, $\hat{y}$ denotes

227    the prediction of each response for each array of predictors, and $\bar{y}$ denotes the mean of

228    true responses in the data. Usually, higher training and testing R-squared values imply

229    lower bias and variance in the predictions, respectively. Fig. 2 shows the whole SL

230    process. Once we select the best SL model with the highest prediction accuracy, we

231    can train it again with all 68 data points for the further analysis.



232

233          Fig. 2 Data analysis procedure for each SL model, where ovals and rounded

234               rectangles denote the input and estimated variables, respectively
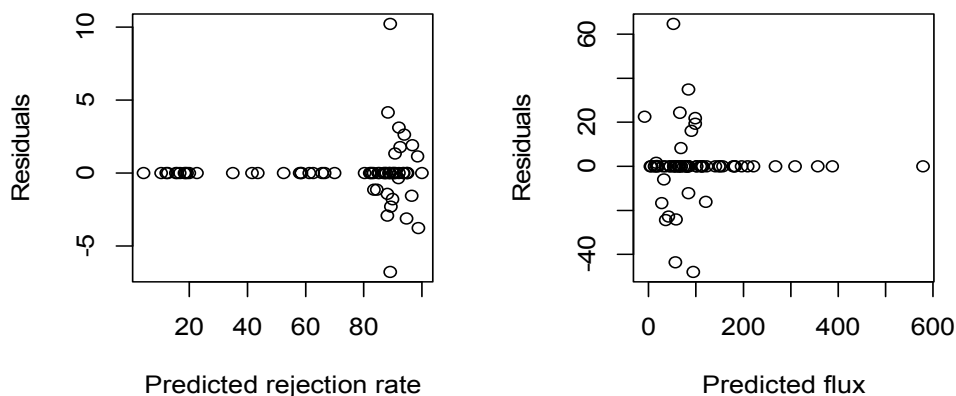
235    **3.2.2 Implementation of supervised learning**

236        Since our data size is small compared to the number of predictors, to avoid over-

237    fitting of NN and SVM, only statistically significant predictors are used for training.

238    Here we apply LR and MART (which are robust to irrelevant predictors) to analyze

239    and extract significant predictors. Also, cross validation is implemented to determine

240    appropriate controlling parameters of NN and SVM, for optimizing the learning

241    performance.

242    **3.2.2.1 Analysis of predictors' significance**

243        According to LR analysis,  the coefficients of PVC Wt% and evaporation time,

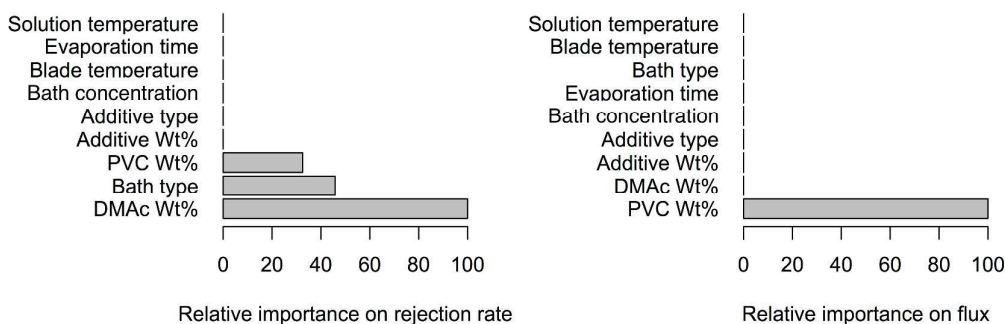244    those of DMAc Wt% and Additive Wt%, and those of additive and bath types are

245　statistically significant at level 0, level 0.01, and level 0.05 for the rejection rate.

246　However, there are only two statistically significant coefficients: one of PVC Wt% at

247　level 0 and another of DMAc Wt% at level 0.1 for the flux. In other words, only a few

248　processing parameters can provide significant information on the predictions;

249　especially for the flux, PVC Wt% and DMAc Wt% are the two that carry the most

250　amount of information. The low statistical significances are partially due to the small

251　number of measurements. The linearity assumption on the relationship can be tested

252　with R-squared values, which we will discuss later. In addition, the identical and

253　independent distribution assumption on the noise can be tested by residual versus

254　predicted response plots, which are shown in Fig. 3. Although the mean of residuals is

255　indeed zero, the variance does not follow the null plot; this may be because our data is

256　collected via a controlled parameter method.



257

258　　　　　　　Fig. 3 Residuals versus Predicted values plots for rejection rate and flux



259

260　　　　　　　　　Fig. 4 Importance plots of predictors on each indicator
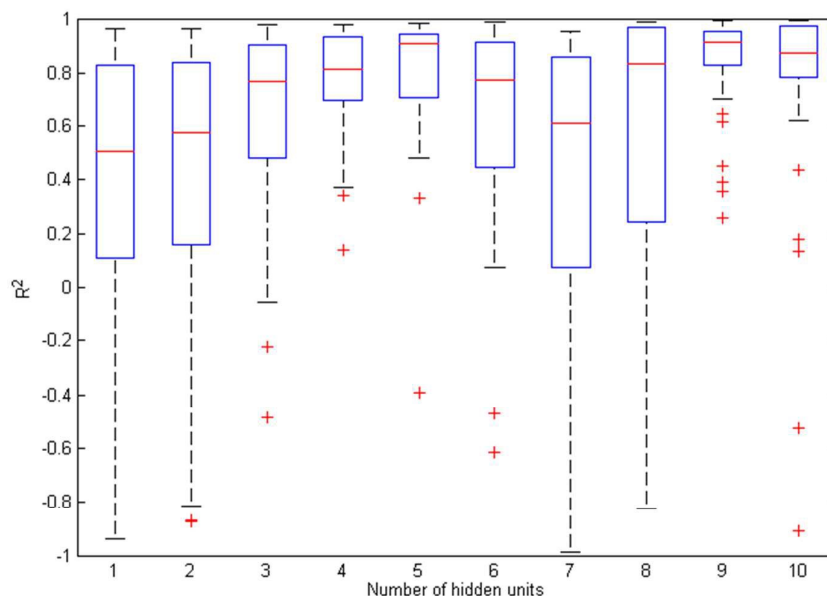
261　　　In case of MART analysis, the resulting importance rankings of each predictor

262　for predictions are shown in Fig. 4. We can see that the number of significant

263　predictors is even fewer than that in LR for each indicator. The importance order is

264   DMAc Wt% > Bath type > PVC Wt% for rejection rate, and only PVC Wt%

265   determines the regression tree for flux.

266       In sumary, LR suggests that PVC Wt% and DMAc Wt% are the two most

267   significant predictors. MART claims that the importance order of predictors is DMAc

268   Wt% > Bath type > PVC Wt% for rejection rate, while only PVC Wt% determines the

269   regression tree for flux. Based on the results of LR and MART, we remove the

270   insignificant predictors (solution temperature) and then train NN and SVM with the

271   appropriate controlling parameters determined by cross validation.

272   **3.2.2.2 Selection of appropriate controlling parameters for NN and SVM**

273       As shown in Fig.1, the responses in our data are correlated, so NN is more

274   appropriate than any other SL model, which can only predict the rejection rate and the

275   flux separately. To apply NN, we should first assume that the categorical predictors

276   (additive type and bath type) are numerical. In addition, we remove the unimportant

277   predicator (solution temperature) and normalize all input predictors to zero-mean and

278   one-standard-deviation.



279

280    Fig. 5 Box-plots of testing R-squared values over 50 training processes with different

281                                    hidden layer sizes

282   Furthermore, we select appropriate controlling parameters. Usually, one hidden layer

283   is sufficient for a small training set. To select the optimal number of hidden units, we

284   repeat the learning processes 50 times for each, and then select the one with a high

285   mean and a low variance of testing R-squared values. During each process, we

286   randomly split the data into a training set, a validation set, and a testing set, with the

287 ratio 51/10/7, and then select the best number of epochs through cross-validation. The
288 resulting box-plots are shown in Fig. 5. We can see the optimal number of hidden
289 units is 9, with both the highest mean (0.8218) and the lowest variance of testing R-
290 squared values.

291     As regard to SVM, since our data size is small, we select only the statistically
292 significant 6 predictors in LR and MART to avoid overfitting. Furthermore, we
293 choose the appropriate controlling parameters with five-fold cross-validation. The
294 resulting support vectors are from all measurements except the $43^{rd}$ or $18^{th}$
295 measurements for the rejection rate or the flux, implying the risk of over-fitting.

296

297 **4. Results and discussion**

298 **4.1 Performance of SL models and selections**

299     The training and testing R-squared values of all SL models introduced above are
300 listed in Table 1, where Rm and Rn denote the training and testing R-squared values,
301 respectively, and y1 and y2 denote the rejection rate and the flux. We can see NN is
302 the best SL model, with the highest Rm and Rn for both y1 and y2. The second best
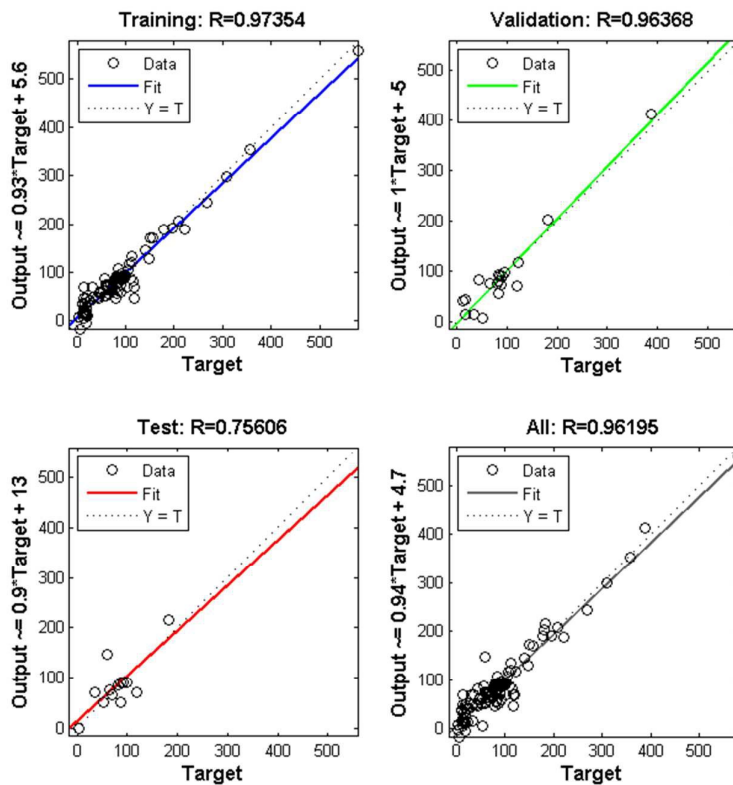303 SL model is SVM, which performs considerably worse for y2 and Rn.

304            Table 1 Summary of performance of different SL models

|           | MART    | NN     | LR     | SVM    |
|-----------|---------|--------|--------|--------|
| **Rm(y1)** | 0.2122  | 0.8897 | 0.6577 | 0.8065 |
| **Rm(y2)** | 0.0725  | 0.9175 | 0.6887 | 0.6583 |
| **Rn(y1)** | 0.0784  | 0.6344 | 0.3104 | 0.4344 |
| **Rn(y2)** | -0.0329 | 0.8093 | 0.1800 | 0.6583 |

305
306     By combining the performance results in Table 1 and the properties of each SL
307 model, we can reveal some interesting underlying characteristics of the data. We
308 begin with the worst SL model, MART, which has very low R-squared values for all
309 conditions. In other words, the piecewise constant approximation does not work on
310 this data, partially due to the small number of controlled measurements. However, we
311 find that both the bias and variance are lower for the rejection rate. Thus, compared to
312 the flux, the rejection rate has relatively high order interactions with processing
313 parameters. This argument can be verified with the performance of LR. Both training
314 R-squared values are relatively high. Especially for the flux, this value is even higher

11

315    than that of SVM. Furthermore, SVM has much higher training R-squared of the

316    rejection rate, and testing R-squared of both rejection and flux than those of LR.

317    Therefore, the relationship between the flux and the processing parameters is

318    approximately linear, but the rejection rate may have more complex and higher order

319    interactions between the processing parameters. In addition, the noise of the

320    measurement data is relatively high. Finally, although the testing R-squared values of

321    SVM are much higher than LR due to the noise reduction in the higher dimensional

322    feature space, they are still much lower than those of NN. This verifies the overfitting

323    of SVM on small data, even when the regularization cost is set as high as $2^5$.
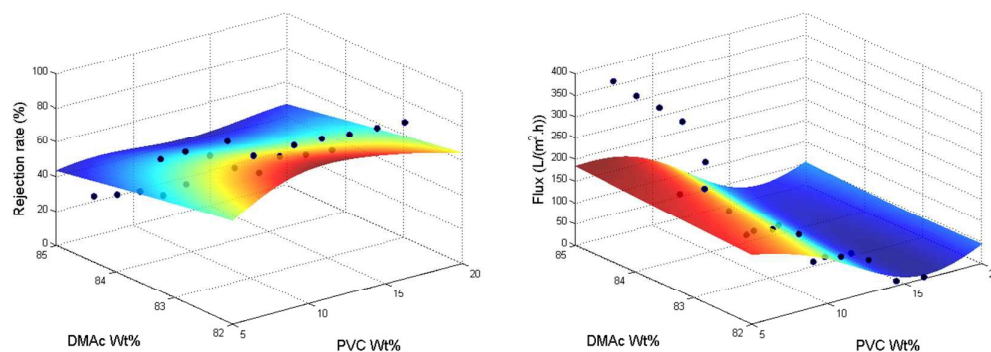


324

325    Fig. 6 Prediction versus response plots for training, validation, testing, and the whole

326    data set; target and output denote the true response and the predicted response by NN,

327                                  respectively

328       NN beats all other SL models in all aspects, and if the whole data is used for

329    training, it has training R-squared values as high as 0.8992 and 0.9559 for the

330    rejection rate and the flux. Thus, compared to the numerical approximation on

331    categorical predictors, the correlation between the rejection rate and the flux is much

332    more important in our predictions. To visualize the performance of NN, we plot the

333    prediction versus the true response in Fig. 6. The performance is considered perfect if

334    the point lies on the line with intersection 0 and slope 1. Furthermore, we plot the

335    training data points and fitting curves of SVM and NN inside the predictor subspace

336    of PVC Wt% and DMAc Wt% in Fig. 7 and Fig. 8 by fixing all other predictors as

337    Additive Wt% = 0%, Additive type = None, Evaporation time = 5 sec, Blade

338    temperature = 60 °C, Bath type = Water, and volume concentration of solute in

339    gelation bath= 0 mg/L. We can see that the fitting curves of NN are smoother and fit

340    the training data better. In summary, because our data set is very small and noisy, the

341    complex relationship between the rejection rate and the processing parameters is hard

342    to fit with a good trade-off between bias and variance. Fortunately, we have the

343    helpful information that tells us that it is correlated with the flux, which has a much

344    simpler linear relationship, so we can apply NN to fit these two indicators.
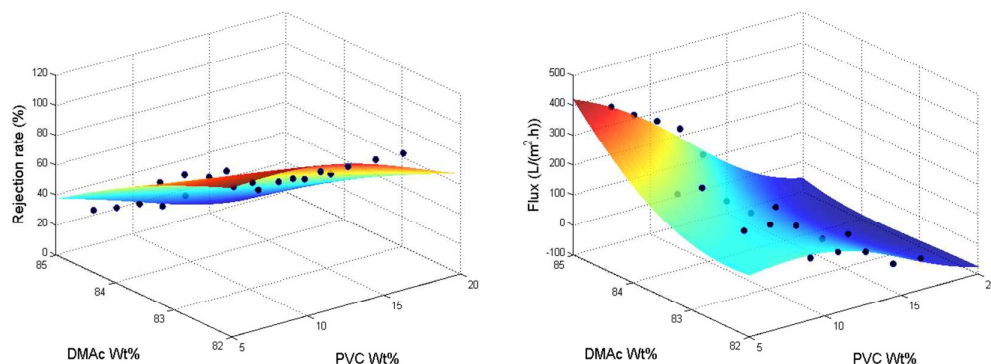
345



346

347        Fig. 7 Training data and fitting curves of rejection rate and flux in the subspace of

348                        PVC Wt% and DMAc Wt% using SVM
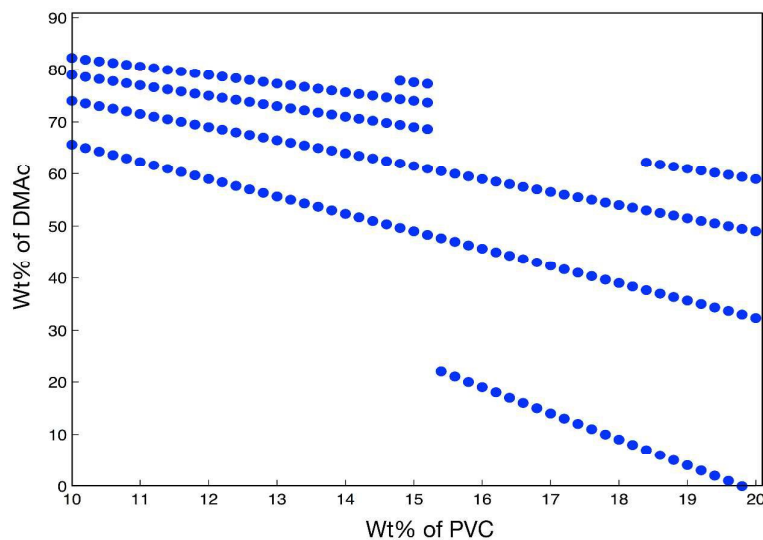
349

350



351

352        Fig. 8 Training data and fitting curves of rejection rate and flux in the subspace of

353                        PVC Wt% and DMAc Wt% using NN

354  **4.2 Optimization with NN**

355  In this section, we use NN model to find the optimal combinations of processing

356  parameters to maximize the flux under the constraint that the rejection rate of BSA

357  should be no less than 80%. The idea is very simple: we search over the predictor

358  space to find certain combinations that achieve the maximum predicted flux under the

359  constraint regarding the predicted rejection rate by NN. For example, when we fix

360  Additive Wt% = 1%, Additive type = PEG600, Evaporation time = 35 sec, Blade

361  temperature = 70 °C, Bath type = Water, and Bath concentration= 0 mg/L, the

362  possible combinations of PVC Wt% and DMAc Wt% satisfying rejection rate >=

363  80%, flux >= 200 L/(m$^2$·h) are scattered in Fig. 9. It is noticed that the combinations

364  are almost impossible in reality in the case of DMAc Wt%<40% or DMAc

365  Wt%>85%. Therefore, a question is raised here on how to perform an efficient and

366  reliable search.  As a matter of fact, regarding to the problem, there exist two main

367  difficulties: (1) when searching over a high-dimensional predictor space, the

368  computation cost is very high; and (2) the predictions have high variance since the

369  size of the training data is small. To overcome these difficulties, we first narrow down

370  the search space by utilizing additional knowledge about the experiments and

371  constraints on predictors. There are several obvious constraints, such as if Additive

372  type = None, then Additive Wt% = 0%; if Bath type = water, then Bath

373  concentration= 0 mg/L. In addition, our focus is on estimating how the addition of

374  PVB into PVC improves the performance of membranes, so we introduce $k$ as the

375  ratio of PVC Wt%/ Polymer Wt%, giving us $0 < k < 1$. Furthermore, we should keep

376  the Polymer Wt% at no greater than 21%. Note that DMAc Wt% can be easily

377  calculated using Eq.3 and Eq.4.

378      So we can use $k$ instead of DMAc Wt%. On the other hand, although the

379  prediction accuracy is not guaranteed over the whole predictor space, both training

380  and testing R-squared are very high within the data set. This means that if the search

381  points are not too far away from the measurement points, the corresponding

382  predictions are reliable. In particular, we have the search space PVC Wt% =

383  7.5:0.5:18 (%), k = $(\lceil$PVC Wt%/21$\rceil)$, 0.05:0.9 , and Additive Wt% = 1:1:5 (%) if

384  Additive type is not None, Evaporation time = 5:15:110 (sec), Blade temperature =

385  30:10:80 (°C), and Bath concentration = 10:10:80 (mg/L).
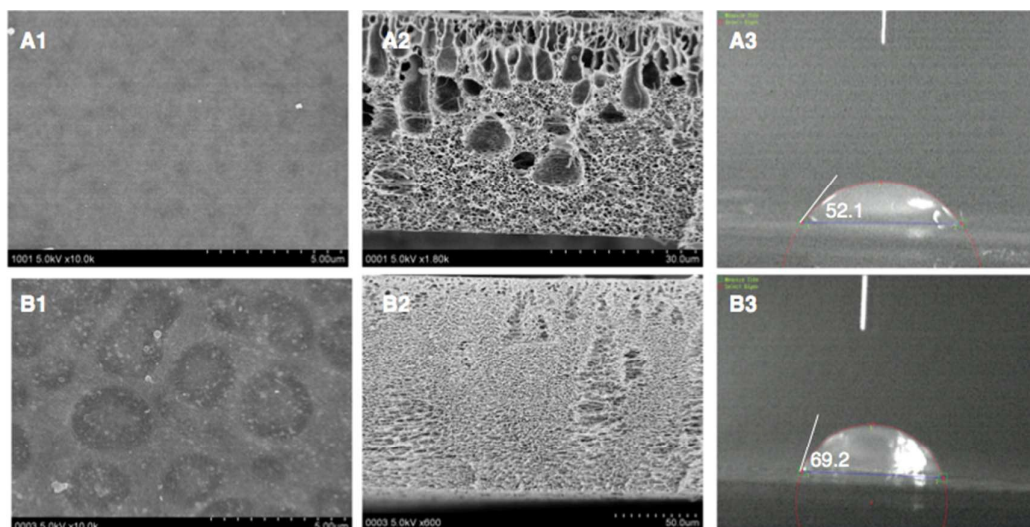
14

386

387    Fig. 9 Possible combinations of PVC Wt% and DMAc Wt% for specific constraints

388                          on indicators fixing all other processing parameters

389          Finally, we select the combination of processing parameters that have the

390    maximum flux under the constraint 80% $\leqslant$ rejection rate $\leqslant$100%. We find with the

391    water bath that the optimal combination of processing parameters is PVC Wt% =

392    7.5%, DMAc Wt% = 84%, Additive Wt% = 1%, k = 0.5 (PVB Wt%=7.5%), Additive

393    type = PEG600, Evaporation time = 5 (sec), and Blade temperature = 30 (°C), leading

394    to the rejection rate = 80.03% and the flux = 329.88 (L/(m$^2$·h)). Similarly, in the

395    DMAc bath, we find that when PVC Wt% = 16%, DMAc Wt% = 78%, Additive Wt%

396    = 2%, k = 0.8 (PVB Wt%=4%), Additive type = PVP k90, Evaporation time = 5 (sec),

397    Blade temperature = 30 (°C), and Bath concentration= 80 (mg/L), we have the

398    rejection rate = 81.39% and the maximum flux = 271.61 L/(m$^2$·h). Although our

399    results are not guaranteed to be globally optimal, they are much robust than the best

400    measurement, which has the rejection rate = 82.07% and the flux = 122.70 L/(m$^2$·h)

401    (with the processing parameters PVC Wt% = 12.6%, DMAc Wt% = 77%, Additive

402    Wt% = 5%, k = 0.7 (PVB Wt%=5.4%), Additive type = PEG600, Evaporation time =

403    10 sec, Blade temperature = 60 °C, Bath type = DMAc, and Bath concentration= 80

404    mg/L). To check the accuracy of the models used to optimize membrane performance,

405    we fabricated PVC/PVB flat sheet membranes strictly under the above optimized

406    parameters. Fig.10 shows the surface and cross-section morphology and the contact

407    angle of the as-prepared membranes. In the case of pure water gelation bath, the

408    rejection rate of the as-prepared membrane was 80.2% and the flux was 318.27

409 L/(m²·h), while in the case of DMAc as the solute of gelation bath, the as-prepared
410 membrane has the rejection rate of 86.2% and the flux of 298.5 L/(m²·h). The results
411 showed that there was a very good agreement between the model predictions and
412 experimental data.

413


414 Fig.10 Morphology and the contact angle of PVC/PVB composite  membranes
415 (A: the membrane prepared under optimized parameters in the case of using pure
416 water as gelation bath, B:  the membrane prepared under optimized parameters in the
417 case of using DMAc as the solute of gelation bath. 1. Suface structure  2. Cross-
418 section structure     3. Contact angle)

419 **5. Conclusions**

420 In this paper, we provide a systematical approach, namely, an SL-based
421 framework for experimental data analytics, for modeling and optimizing membrane
422 responses for complex combinations of membrane features during fabrication. This
423 approach consists of the following procedures. First, control experiments are
424 established to get various membranes with differing performances by combining
425 various fabrication conditions. Second, the characteristics of the feature variables are
426 analyzed in order to ascertain the quality of the data, as well as the data dependencies
427 among the variables. Third, four SL approaches (MART, NN, LR, SVM) are
428 employed to systematically analyze membrane performance and fabrication
429 conditions in a rigorous fashion. Finally, the most reliable and trustful SL model is
430 selected to optimize the fabrication conditions and predict the most favorable
431 performance of PVC/PVB ultrafiltration membranes. During this last step, we analyze
432 multiple responses simultaneously with multiple input feature variables. In this way,

433　we eliminate most unnecessary assumptions that are traditionally proposed by other

434　methods. In addition, this approach simplifies the analysis process by using a unified

435　SL framework that has been thoroughly investigated by machine learning

436　communities[24]. This advantage surpasses previously reported DOE approaches in that

437　these standard SL approaches provide smaller biases and variances for data analysis.

438　Thus, the SL approaches offer us a more standard method not only in procedure but

439　also with more rigorous results.

440　　　Additionally, we glean several interesting findings from this research. One is

441　how to find the optimal mixture of feature compounds for the fabrication processes

442　more effectively and efficiently. Another is that among the tested SL approaches, the

443　NN method provides the most reliable and trusted results. In the future, we will

444　investigate how to develop a recursive and automated data-driven experimental

445　analytics approach to design performance-specific membranes more effectively and

446　efficiently.

447

457

458　**References:**

459　1　S. Ramesh, A. H. Yahaya and A. K. Arof, *Solid State Ionics*, 2002, **148**, 483–486.

460　2　Z. Yu, X. Liu, F. Zhao, X. Liang and Y. Tian, *J. Appl. Polym. Sci.*, 2015, **132**,
461　　41267.

462　3　N. Wang, A. Raza, Y. Si, J. Yu, G. Sun and B. Ding, *Journal of Colloid and*
463　　*Interface Science*, 2013, **398**, 240–246.

464　4　S. Chuayjuljit, R. Thongraar and O. Saravari, *Journal of Reinforced Plastics and*
465　　*Composites*, 2008, **27**, 431–442.

466　5　M. Jakic, N. S. Vrandecic and I. Klaric, *Polymer Degradation and Stability*,
467　　2013, **98**, 1738–1743.

468　6　X. Zhao and N. Zhang, *Journal of Tianjing University of Science and*
469　　*Technology*, 2007, **22**, 36–39.

470    7    J. zhu, L. Chi, Y. Zhang, A. Saddat and Z. Zhang, *Water Purification*
471         *Technology*, 2012, **31**, 46–54.
472    8    Y. Peng and Y. Sui, *Desalination*, 2006, **196**, 13–21.
473    9    Y. Sui, *Beijing University of Technology Thesis for Master Degree*, Beijing,
474         China, 2004, 24–26.
475   10    X. Zhao and K. Xu, *Plastics Sci.& Technology*, 2010, 1–6.
476   11    E. Corradini, A. F. Rubira and E. C. Muniz, *European polymer journal*, 1997, **33**,
477         1651–1658.
478   12    S.Y. L. Leung, W.H. Chan and C.H. Luk, *Chemometrics and Intelligent*
479         *Laboratory Systems*, 2000, **53**, 21–35.
480   13    S.Y. Lam Leung, W.H. Chan, C.H. Leung and C.H. Luk, *Chemometrics and*
481         *Intelligent Laboratory Systems*, 1998, **40**, 203–213.
482   14    W.H. Chan and S.C. Tsao, *Chemometrics and Intelligent Laboratory Systems*,
483         2003, **65**, 241–256.
484   15    M. Khayet, C. Cojocaru, M. Essalhi, M. C. García-Payo, P. Arribas and L.
485         García-Fernández, *DES*, 2012, **287**, 146–158.
486   16    L. Wenjau and O. Soonchuan, *Computer and Automation Engineering (ICCAE),*
487         *2010 The 2nd International Conference*, 2010, **2**, 50–54.
488   17    P. W. Araujo and R. G. Brereton, *Trends in Analytical Chemistry*, 1996, **15**, 63–
489         70.
490   18    K. I. Wong, P. K. Wong, C. S. Cheung and C. M. Vong, *Energy*, 2013, **55**, 519–
491         528.
492   19    Y. Reich and S. V. Barai, *Artificial Intelligence in Engineering*, 1999, **13**, 257–
493         272.
494   20    B. L. Whitehall, S.-Y. Lu and R. E. Stepp, *Artificial Intelligence in Engineering*,
495         1990, **5**, 189–198.
496   21    F. J. Alexander and T. Lookman, *Novel Approaches to StatisticalLearning in*
497         *Materials Science*, Informatics for Materials Science and Engineering, 2013.
498   22    S. S. Madaeni, N. T. Hasankiadeh, A. R. Kurdian and A. Rahimpour, *Separation*
499         *and Purification Technology*, 2010, **76**, 33–43.
500   23    X. Xi, Z. Wang, J. Zhang, Y. Zhou, N. Chen, L. Shi, D. Wenyue, L. Cheng and
501         W. Yang, *DESALINATION AND WATER TREATMENT*, 2013, **51**, 3970–3978.
502   24    C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer Verlag,
503         2006.
504

## Figures

Fig. 1 Scatter plots over measurements

Fig. 2 Data mining procedure for each SL model, where ovals and rounded rectangles denote the input and estimated variables, respectively

Fig. 3 Residuals versus Predicted values plots for rejection rate and flux

Fig. 4 Importance plots of predictors on each indicator

Fig. 5 Box-plots of testing R-squared values over 50 training processes with different hidden layer sizes

Fig. 6 Prediction versus response plots for training, validation, testing, and the whole data set; target and output denote the true response and the predicted response by NN, respectively

Fig. 7 Training data and fitting curves of rejection rate and flux in the subspace of PVC Wt% and DMAc Wt% using SVM

Fig. 8 Training data and fitting curves of rejection rate and flux in the subspace of PVC Wt% and DMAc Wt% using NN

Fig. 9 Feasible combinations of PVC Wt% and DMAc Wt% for specific constraints on indicators fixing all other processing parameters

Fig. 10 Morphology and the contact angle of the as-prepared optimized membranes (A: the membrane prepared under optimized parameters in the case of using pure water as gelation bath, B: the membrane prepared under optimized parameters in the case of using DMAc as the solute of gelation bath. 1. Suface structure, 2. Cross-section structure, 3. Contact angle.)
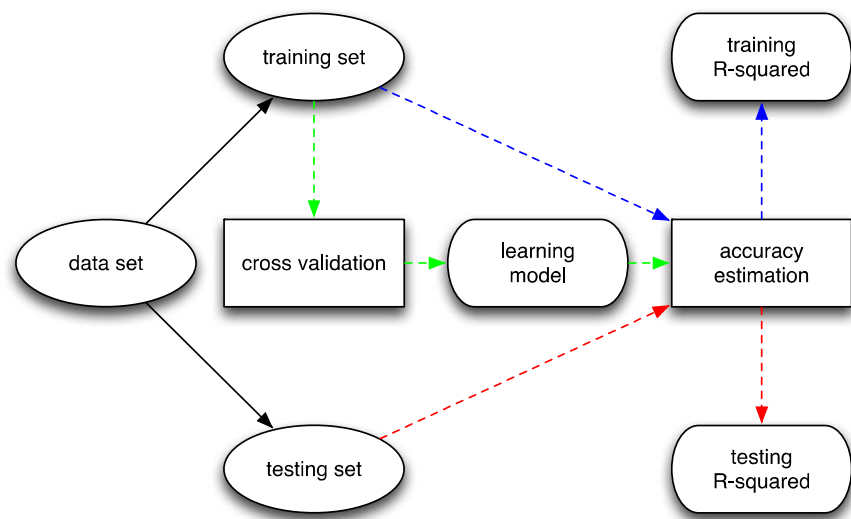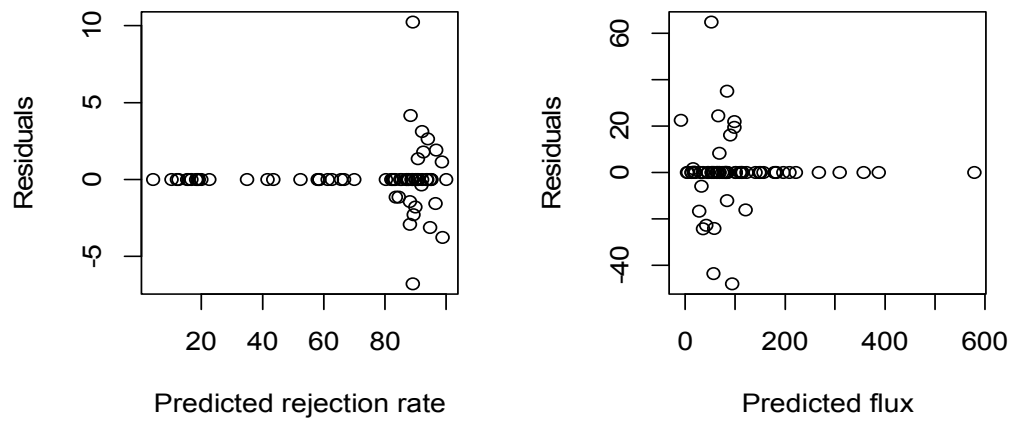
Figure 1

Figure 2

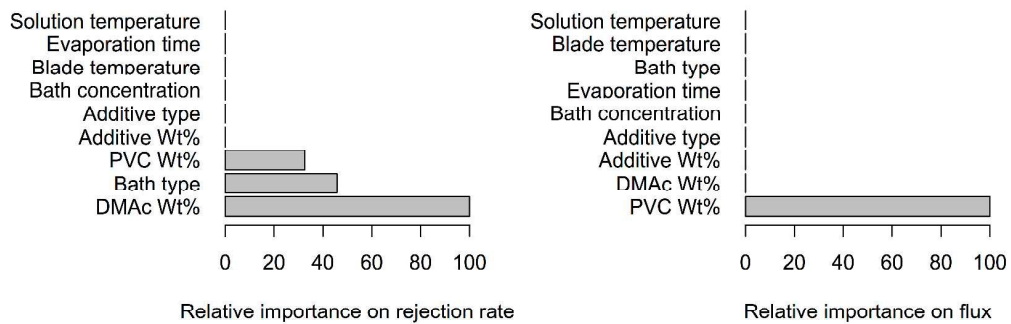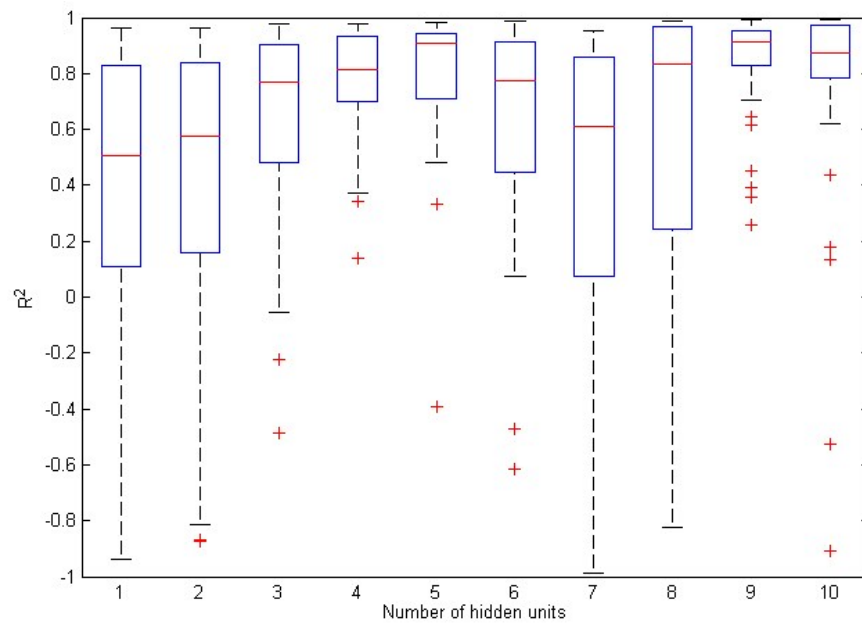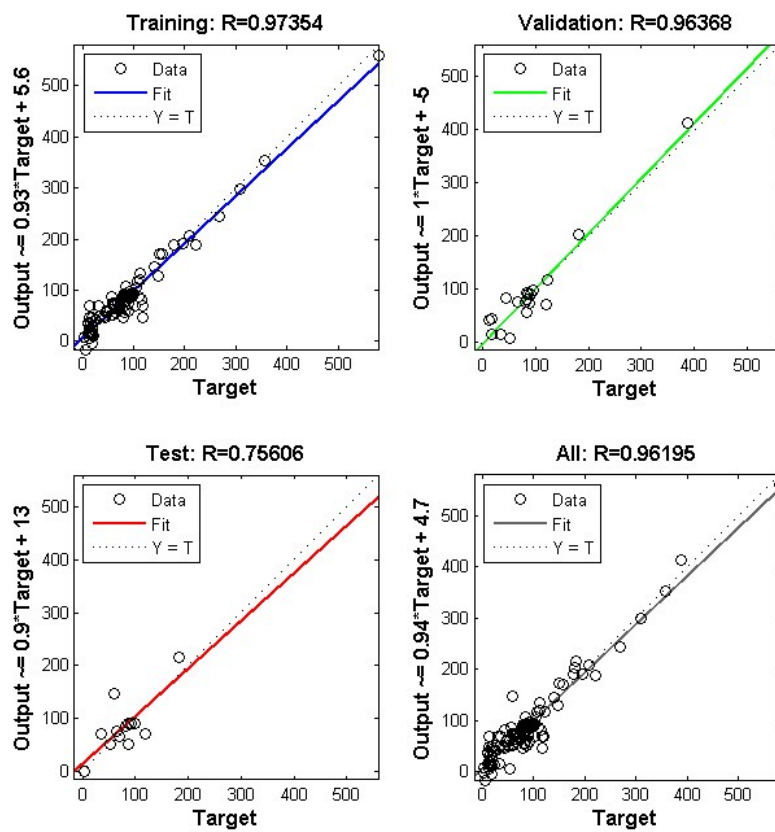Figure 3



Predicted rejection rate                    Predicted flux
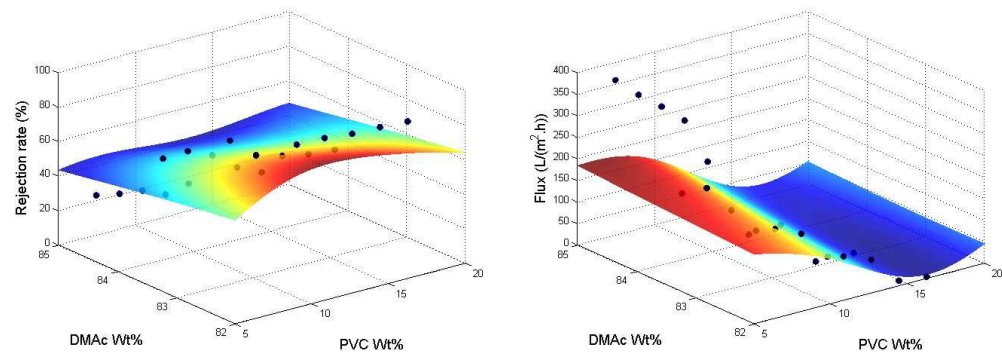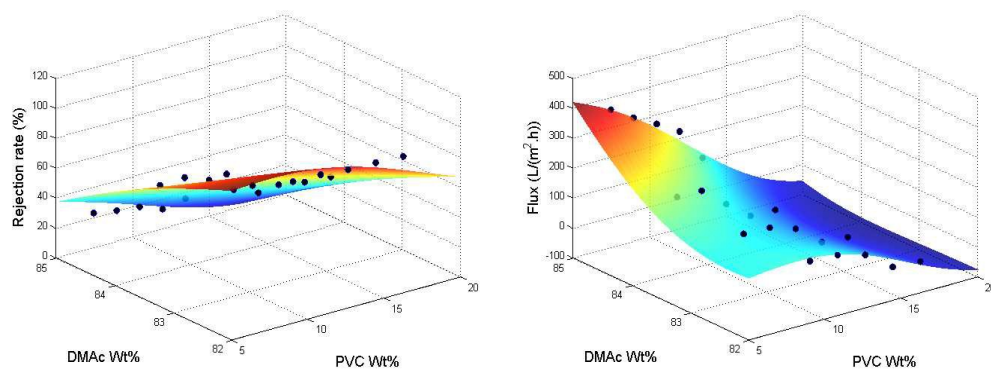
Figure 4

Figure 5

Figure 6

Figure 7

Figure 8

Figure 9

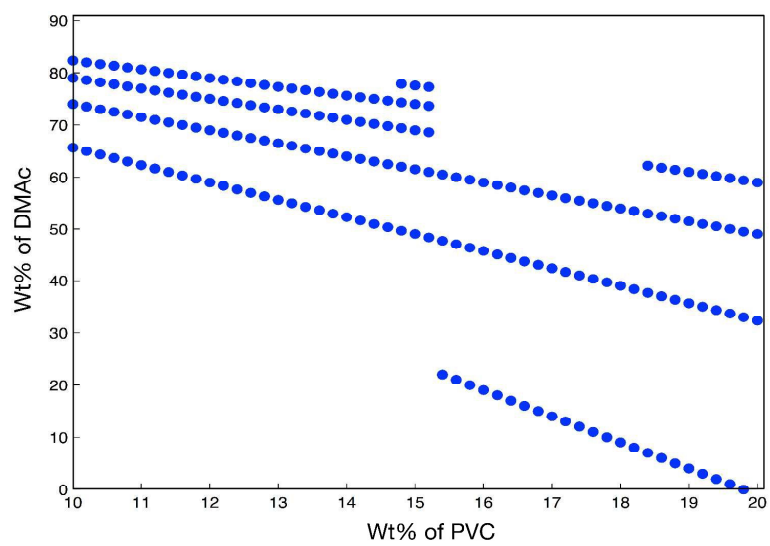Fig. 10

**Tables**

Table.1 Summary of performance of different SL models

Table 1

|          | MART    | NN     | LR     | SVM    |
|----------|---------|--------|--------|--------|
| **Rm(y1)** | 0.2122  | 0.8897 | 0.6577 | 0.8065 |
| **Rm(y2)** | 0.0725  | 0.9175 | 0.6887 | 0.6583 |
| **Rn(y1)** | 0.0784  | 0.6344 | 0.3104 | 0.4344 |
| **Rn(y2)** | -0.0329 | 0.8093 | 0.1800 | 0.6583 |

**Graphical abstract**