# Journal Pre-proof
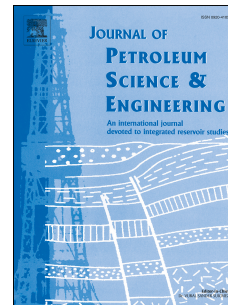
Prediction of CO$_2$ diffusivity in brine using white-box machine learning

Menad Nait Amar, Ashkan Jahanbani Ghahfarokhi

Please cite this article as: Amar, M.N., Ghahfarokhi, A.J., Prediction of CO$_2$ diffusivity in brine using white-box machine learning, *Journal of Petroleum Science and Engineering* (2020), doi: https://doi.org/10.1016/j.petrol.2020.107037.

# Prediction of CO$_2$ Diffusivity in Brine Using White-Box Machine Learning

Menad Nait Amar [a] and Ashkan Jahanbani Ghahfarokhi [b][*]

*a Département Etudes Thermodynamiques, Division Laboratoires, Sonatrach, Avenue 1er Novembre, 35000,Boumerdes, Algeria*

*b Department of Geoscience and Petroleum, Norwegian University of Science and Technology (NTNU), S. P. Andersens veg 15b, 7031 Trondheim, Norway*

**Abstract**

Accurate knowledge of the diffusivity coefficient of CO$_2$ in brine has a significant effect on the design success and monitoring of CO$_2$ storage in saline aquifers, which is a part of carbon capture and sequestration (CCS). Frequently applied experimental approaches for determining this parameter are expensive and time-consuming, and empirical models cannot ensure accurate predictions. Therefore, there is a need to establish cutting-edge correlations for prediction of the diffusivity coefficient of CO$_2$ in brine under various operating conditions. In this work, two white-box machine learning techniques, namely group method of data handling (GMDH) and gene expression programming (GEP) were implemented for correlating the diffusivity coefficient of CO$_2$ in brine with pressure, temperature and the viscosity of the solvent. The obtained results demonstrated the accuracy of the proposed correlations. In addition, statistical and graphical analysis of the performances revealed that GEP correlation outperforms the GMDH correlation, decision trees (DTs), random forest (RF) and all the previous predictive models. GEP correlation exhibited an overall average absolute relative deviation (AARD) of 4.3014% and coefficient of determination (R$^2$) of 0.9979. Finally, by performing the outliers detection, the validity of the GEP correlation was confirmed and only two experimental data points were identified as outliers.

**Keywords** – CO$_2$-brine; diffusivity coefficient; empirical correlations; GEP; GMDH.

*Corresponding author: ashkan.jahanbani@ntnu.no

## 1. Introduction

41

42    The increased amount of $CO_2$ in the atmosphere is considered the most important

43    concern of the 21$^{st}$ century around the globe (Boot-Handford et al., 2014). This issue results

44    mainly from fossil fuels being used continuously in different sectors of industry (Azzolina et

45    al., 2015; Boot-Handford et al., 2014; Nait Amar et al., 2019a; Venkatraman and Alsberg,

46    2017). Therefore, some new approaches have emerged as promising means for reducing the

47    $CO_2$ levels in the atmosphere, among which, carbon capture and sequestration (CCS) in

48    geological formations is still the most attractive one (Amini et al., 2012; Grude et al., 2014;

49    Shahkarami et al., 2014).

50    CCS has gained much interest within many fields of engineering, environment and

51    energy (Bhakta et al., 2015; Davarazar et al., 2019; Gibbins and Chalmers, 2008; Lee et al.,

52    2010; Mohagheghian et al., 2019; Riahi et al., 2004). The sequestration of the captured $CO_2$ in

53    saline aquifers is the most frequently applied strategy while implementing CCS (Gershenzon

54    et al., 2015, 2014).  In addition, some other applications of $CO_2$ such as enhanced oil recovery

55    methods based on $CO_2$ injection, have served as useful and vital means in reduction of $CO_2$

56    levels in the atmosphere (Bachu et al., 2004; Ettehadtavakkol et al., 2014; Gozalpour et al.,

57    2005; Holtz et al., 2001; Nait Amar and Zeraibi, 2019).

58    The process of $CO_2$ injection in saline aquifers is subjected to some mechanisms,

59    namely the contact of $CO_2$ with the in-situ water and its dissolution in brine through

60    molecular diffusion. Therefore, accurate knowledge and determination of the parameters that

61    play important role while monitoring CCS are of vital interest. Diffusivity coefficient which

62    characterizes the diffusivity of fluid is one of these parameters (Cadogan et al., 2014a;

63    Guzmán and Garrido, 2012). Indeed, this parameter has a noticeable effect on the chemical

64    reactions and the interfacial mass transfer occurring deep underground, in addition to

65    impacting the flow path, transport behavior and the quantitative description of diffusion

66  during $CO_2$ injection (Farajzadeh et al., 2009; Guzmán and Garrido, 2012; Mutoru et al.,

67  2011; Trevisan et al., 2014).

68  Determination of diffusivity coefficient of the $CO_2$ in brine can be done by two

69  distinguished approaches. The first one consists of performing experimental measurements

70  and the second approach is by means of the empirical models. The lab measurements are

71  divided into direct and indirect tests (Lu et al., 2013). The direct procedures are based on the

72  accurate knowledge of the gas concentration in the solvent (Cadogan et al., 2014b; Frank et

73  al., 1996), while the indirect tests exploit the gained information related to the diffusivity of

74  gas such as interfacial tension, volumes of gas and liquids, and the operational conditions

75  (Jang et al., 2018; Li et al., 2016). The experimental approaches are known to deliver accurate

76  results. However, their implementation is time-consuming and demands sophisticated

77  equipment. As a result, several researchers have developed empirical models for estimating

78  the diffusivity coefficient of CO2 in brine. Wilke and Chang (1955) established a predictive

79  correlation for estimating the diffusivity coefficient in numerous dilute solutions. The model

80  employs viscosity and temperature as input parameters.  Lu et al. (2013) developed a model

81  for estimating the diffusivity coefficient of $CO_2$ in water without considering the pressure

82  effect. There model is applicable to cases with temperatures between 268K and 473K. An

83  extended version of  Lu et al. (2013) model at high pressure and temperature conditions, was

84  proposed by Moultos et al. (2016) by applying the concept of molecular dynamics

85  simulations. Although the correlations developed by Lu et al. (2013) and Moultos et al. (2016)

86  are accurate for the $CO_2$-pure water system, they cannot be applied to cases where brine is the

87  solvent. Cadogan et al. (2014a) utilized experimental results for $CO_2$ diffusivity coefficients

88  against brine viscosity at temperature of 298K for establishing a modified Stokes-Einstein

89  relation. Table 1 reports the mathematical formulations of the above models for predicting the

90  diffusivity coefficient of $CO_2$ in water and brine.

91

92    An in-depth review of the available correlations for predicting the diffusivity coefficient

93    of $CO_2$ in brine reveals the limitations of these techniques from the applicability and accuracy

94    perspectives (Feng et al., 2019).

95    In recent years, researchers have shown an increased interest in the application of

96    machine learning techniques for modeling complex systems (Jeong et al., 2018; Nait Amar et

97    al., 2019b; Nait Amar and Zeraibi, 2019; Nomeli and Riaz, 2017; Piotrowski and

98    Napiorkowski, 2012). Machine learning techniques can be divided into computer-aided

99    methods such as support vector regression (SVR) and decision tree, and explicit methods such

100   as gene expression programming (GEP) and group method of data handling (GMDH). The

101   first category is known as black-box approaches, and this means that their paradigms are

102   dependent on a computer-aided technique, while the second category is recognized as white-

103   box methods which means that they deliver explicit expressions (Nait Amar et al., 2019c).

104   Recently, Feng et al. (2019) developed a predictive model for estimating the diffusivity

105   coefficient of $CO_2$ in brine by coupling genetic algorithm with mixed kernel SVR and they

106   obtained satisfactory performances. However, the application and accuracy of their

107   established model depend on calculability efforts, and this presents an issue in terms of

108   flexibility for further utilization.

109   The main contribution and novelty of this study consist of establishing two distinct

110   explicit and simple-to-use correlations for accurate prediction of diffusivity coefficient of $CO_2$

111   in brine under various operational conditions. To do so, group method of data handling

112   (GMDH) and gene expression programming (GEP) were implemented with three input

113   parameters, namely pressure, temperature and viscosity of the solvent, using a representative

114   experimental database. Besides, decision trees (DTs) and random forest (RF) were considered

115   for comparison with the best-result explicit correlation. Statistical and graphical assessment

116   criteria were applied for evaluating the newly proposed correlations and compared their

4

117 performances with prior paradigms. Lastly, Leverage approach was performed to verify the

118 quality of the employed experimental data points and define the realm of application for the

119 best fit established correlation.

120      The rest of the paper includes 4 sections. Section 2 describes the database which was

121 utilized for developing the correlations. Section 3 briefs the two applied white-box machine

122 learning techniques, namely GMDH and GEP, and the procedure of their implementation in

123 our study. Results are given and discussed in Section 4. The paper ends with Section 5 which

124 summarizes the main findings.

## 2. Data collection and preparation

126      In the present work, a representative experimental database was collected from the

127 published literature (Cadogan et al., 2015; Cadogan et al., 2014b; Choudhari and

128 Doraiswamy, 1972; Frank, Marco J W and Swaaij, 1996; Lu et al., 2013; Maharajh, 1975;

129 Maharajh and Walkley, 1972; Nijsing et al., 1959; Reddy and Doraiswamy, 1967; Tamimi et

130 al., 1994; Tan and Thorpe, 1992; Thomas and Adams, 1965; Versteeg and van Swaal, 1988;

131 Vivian and Peaceman, 1956; Yang et al., 2006) to develop accurate explicit correlations for

132 estimating the diffusivity coefficient of $CO_2$ in brine at different operating conditions. The

133 database englobes 92 experimental data points with different operating conditions, namely

134 pressure, temperature and viscosity of the solvent. In the context of the affecting variables,

135 salinity affects the solubility, interfacial tension and phase equilibria, thus influencing the

136 diffusivity. In addition, salinity of the solvents affects brine's viscosities (Cadogan, 2015;

137 Feng et al., 2019); therefore, the salinity effect on diffusivity of $CO_2$ in brine is emulated by

138 considering the brine's viscosities as an input parameter while establishing the correlations

139 and the paradigms. The data points collected from previous experimental studies were

140 obtained using various techniques and equipments such as Taylor dispersion, a modified

141 version of Ringborm's apparatus, laminar jet apparatus, laminar falling film, laser-induced

142  fluorescence (LIF), 13C pulsed-field gradient NMR, physical absorption experiments in a

143  stirred vessel operated with a horizontal gas-liquid interface, optical capillary cell via time-

144  dependent Raman spectroscopy, wetted sphere apparatus, Taylor−Aris dispersion method and

145  see-through windowed high-pressure cell. Table 2 reports a detailed statistical insight about

146  the collected data points. In addition, to provide an insightful description of the database, Fig.

147  1 illustrates frequency histograms of the collected dataset and Fig. 2 demonstrates the

148  correlation between diffusivity coefficient and the considered independent variables through

149  cross plots. According to the histograms shown in Fig. 1, it can be seen that the pressure and

150  temperature are mainly distributed in the medium and low ranges, while viscosity shows a

151  symmetric distribution near its mean. According to Fig. 2 (a), pressure exhibits moderate

152  direct relation with diffusivity coefficient. This analogy is more significant for temperature as

153  can be seen from Fig. 2 (b). From Fig. 2 (c), it can be noted that viscosity of solvent has a

154  negative direct relation with diffusivity coefficient.

## 3. Methodology

### 3.1. Group Method of Data Handling (GMDH)

157  Group Method of Data Handling (GMDH) is one of the artificial neural network (ANN)

158  types, which is known to generate explicit correlation between input and output parameters of

159  a given system.  The resulting correlation by applying a GMDH model takes the form of a

160  polynomial (Dargahi-Zarandi et al., 2017). As an ANN type, GMDH structure involves nodes

161  as basic elements for processing the information. These nodes are arranged in different layers

162  from the input layer to the output layer, with or without intermediate layers (Nait Amar et al.,

163  2019c). The GMDH hybrid version (HGMDH) allowed the improvement of the predictability

164  which was somehow insufficient in the first version developed by (Ivakhnenko, 1971). In

165  HGMDH, the interactions among nodes from different layers are allowed. This procedure

166 brings more robustness when modeling complex cases (Rostami et al., 2019). The

167 mathematical form of HGMDH is expressed as shown below:

168
$$y_i = a + \sum_{i=1}^{d}\sum_{j=1}^{d}...\sum_{k=1}^{d}\vartheta_{ij...k}x_i^m x_j^m ... x_k^m \quad m = 1,2,...,2^p \tag{1}$$

169 where $x$ and $y$ stand for the input and the output parameters, respectively; $\vartheta_{ij...k}$ correspond to

170 the polynomial coefficients; $p$ is the number of layers and $d$ is the number of variables.

171 The following points summarize the calculation procedure in HGMDH with a second order:

172 - The following equation defines the expression of a node $N_i$ covering two inputs:

173
$$N_i^{GMDH} = \alpha_0 + \alpha_1 x_i + \alpha_2 x_j + \alpha_3 x_i x_j + \alpha_4 x_i^2 + \alpha_5 x_j^2 \tag{2}$$

174 - Calculation of polynomial coefficients: least square method is applied to calculate the

175 resulting coefficients in the expressions of the different nodes. The following formula

176 is adapted:

177
$$\Delta_j^2 = \sum_{i=1}^{N}\left(N_i^{GMDH} - y_i\right)^2 \qquad j = 1,2,...,\binom{d}{2} \tag{3}$$

178 where $d$ and $N$ are the number of variables and data points, respectively.

179 - Matrix transformation: in order to achieve the final expression, the above equation is

180 transformed to a matrix form (Dargahi-Zarandi et al., 2017; Hemmati-Sarapardeh and

181 Mohagheghian, 2017):

182
$$Y = A^T X \tag{4}$$

183 - The final solution is obtained as follows:

184
$$A^T = yX^T(XX^T)^{-1} \tag{5}$$

185 where $y = \{y_1, y_2, ..., y_d\}$ and $A = \{\alpha_0, \alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5\}$.

186 **3.2. Gene Expression Programming (GEP)**

7

187    Gene expression programming (GEP) is a prevalent evolutionary technique which aims

188    at modeling the systems with accurate explicit expressions. GEP, introduced by Ferreira

189    (Ferreira, 2001), can be regarded as the improved version of genetic programming (GP) which

190    was proposed by Koza (Koza, 1992). GEP employs the standard genetic operators, namely

191    selection, crossover, elitism and mutation, in addition to new implemented actions such as

192    insertion and transposition to search for the reliable correlations. Besides, the codification of

193    the individuals in GEP is performed in the form of chromosomes. In addition, Expression

194    Tree (ET) is introduced for transforming the individuals to real expressions. It is worth

195    mentioning that the genes have a fixed length with terminals which show the variables, such

196    as $\{x_1, x_2, x_3\}$, and operators such as $\{+, /, \times, -, \sqrt{.}, ln\}$ (Teodorescu and Sherwood, 2008).

197    The main steps of GEP are briefed below:

198    -    Initial population: an initial population of correlations codified in the form of

199         chromosomes is generated randomly. The prediction quality of the correlations differs

200         according to a fitness function.

201    -    Standard operators: the well-known genetic operator including elitism, selection,

202         crossover and mutation are applied. Elitism consists of safeguarding the best

203         correlations for the next generations. Selection allows identification of the correlations

204         to be subjected to reproduction for giving new correlations. Crossover is summarized in

205         the process of exchanging parts between two or more correlations, while mutation is

206         done by means of modifying parts of correlations.

207    -    Transposition and insertion: it consists of jumping and activating parts of the genome in

208         the chromosome (Ferreira, 2001).

209    The resulting population is assessed and the operators are reiterated until satisfying a

210    stopping condition.

211

## 3.3. Implementation procedure

212

213    As previously highlighted, the aim behind applying GMDH and GEP as white-box

214    machine learning method is to develop an explicit correlation for accurate prediction of $CO_2$

215    diffusivity in brine under various conditions including pressure, temperature and viscosity.

216    Therefore, the following form is admitted for the two correlations:

217    $$D_c = f(P, T, \mu) \tag{6}$$

218    In the above equation, $D_c$ points out the $CO_2$ diffusivity coefficient in brine expressed in $m^2/s$,

219    and $P, T$ and $\mu$ represent the input parameters of the correlations, viz. pressure, temperature

220    and viscosity, respectively. The input parameters are expressed in MPa, K and mPa.s,

221    respectively.

222    The collected experimental data points that describe these conditions and the obtained

223    diffusivity coefficient of $CO_2$ in brine were prepared for the development of these

224    correlations. The database was divided randomly into a training set with 80% of the whole

225    experimental data points and a testing set which covers the remaining 20%. This dataset

226    partitioning exhibits usually very satisfactory results (Aminu et al., 2019; Benamara et al.,

227    2019; Dargahi-Zarandi et al., 2017; Hemmati-Sarapardeh et al., 2018; Mirjalili, 2015; Yan et

228    al., 2006). Besides, in order to substantiate the better performance and robustness of applied

229    techniques, sensitive analysis of the latter on database were performed.

230    During the development of the GEP correlation, mean square error (MSE) was the

231    considered fitness function for assessing the chromosomes. MSE is expressed as:

232    $$MSE = \frac{\sum_{i=1}^{n}(t_i - o_i)^2}{n} \tag{7}$$

233    where $t_i$ and $o_i$ stand for the measured and the predicted diffusivity coefficient of the $CO_2$ in

234    brine, respectively, and $n$ represents the number of data points.

9

235  While developing the GMDH-based correlation, the number of inputs that can be

236  introduced in the hidden and output nodes was specified to three, while the best highest order

237  of the model was investigated by performing a sensitivity analysis.

238  As indicated, the control parameters of GEP affect the prediction capability. During the

239  establishment of the GEP-based correlation for the prediction of the diffusivity coefficient of

240  $CO_2$ in brine, it was noticed that an increase in the number of chromosomes in the population,

241  the numbers of genes as well as the maximum depth of ET, influence the run-time, the

242  accuracy and the complexity of the generated correlations. Accordingly, these control

243  parameters were tuned. Table 3 reports the final setting of GEP. Consequently, we applied

244  tree encoding, 100 chromosomes, 12 genes, MSE as fitness functions, a function set including

245  $+, -, \times, /, \exp., \sqrt{\phantom{x}}, INV, ln$, and two point mutation, while the stopping criterion was the

246  maximum number of generations (420).

## 4. Results and discussion

### 4.1. Development and evaluation of the correlations

249  Several statistical indexes including average absolute relative deviation (AARD),

250  coefficient of determination ($R^2$) and root mean square error (RMSE), were considered for

251  assessing the quality of the predictions of the newly proposed correlations. These statistical

252  criteria are defined as follows:

$$AARD\% = \frac{1}{n}\sum_{i=1}^{n}\left|\frac{t_i - o_i}{t_i}\right| \times 100 \qquad (8)$$

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(t_i - o_i)^2}{\sum_{i=1}^{n}(o_i - \bar{t})^2} \qquad (9)$$

$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(t_i - o_i)^2} \qquad (10)$$

10

256    The graphical representations of the outcomes of the correlations were illustrated for a

257    visual evaluation and a comparison of the performances. The graphical assessment of the

258    performances was done by means of cross plots, relative error distribution, cumulative

259    frequency distribution of the absolute relative error and bar plots. Cross plots give an insight

260    about the reliability of the correlations by representing their predictions against the real values

261    of the output, and then comparing the obtained distribution against the line $Y = X$ which

262    shows the perfect paradigm. The diagram of relative error distribution details the repartition

263    of the relative error through the output values. Satisfactory distribution of the relative error

264    around the zero-error line indicates the reliability of the correlations. The cumulative

265    frequency diagram of the absolute relative error allows the identification of portions of the

266    data points which are predicted with some pre-specified values of the absolute relative error.

267    Predicting a great portion of the data points with a low absolute relative error value ensures

268    the high reliability of the correlations. Bar plots summarize the statistical criteria of the

269    correlations in visual comparative schemes.


270    As stated in the previous section, sensitivity analyses were conducted to find the best

271    highest order for GMDH and to investigate GEP and GMDH robustness according to samples

272    considered for training and testing phases. To this end, ten realizations were run at each

273    different GMDH highest order, namely two, three and four. The performance comparison

274    through final overall AARD values is depicted in the box plot of Fig. 3. In this figure, the box

275    exhibits specific interquartile range. The red horizontal line denotes the median value, while

276    the lower and higher black horizontal lines represent the best and worst overall AARD values.

277    As can be seen from this figure, considering three as the highest GMDH order results in the

278    best performance with lowest overall AARD value. Therefore, the highest order of the model

279    was set to three. Fig. 4 shows a bar plot comparing AARD values of GEP and GMDH during

280    training and testing phases in the considered ten runs (with different training and testing

11

281   samples). According to this figure, the best performance of GEP was achieved in the fourth

282   run while GMDH showed its best reliability in the third run. Accordingly, these two best-

283   result models were kept for further comparison.

284       Fig. 5 schematizes the best resulted GMDH model for predicting the diffusivity

285   coefficient of $CO_2$ in brine.

286       According to this figure, no middle nodes were resulted in the model. In the same

287   context, the following equation defines the explicit expression of the GMDH model:

288      $D_c = 10^{-9} \times [C_0 + C_1 \times \mu + C_2 \times T + C_3 \times P + C_4 \times T \times \mu + C_5 \times P \times \mu + C_6 \times P \times T + C_7 \times$

289   $\mu^2 + C_8 \times T^2 + C_9 \times P^2 + C_{10} \times P \times T \times \mu + C_{11} \times T \times \mu^2 + C_{12} \times T^2 \times \mu + C_{13} \times P \times \mu^2 + C_{14} \times$

290   $P \times T^2 + C_{15} \times P^2 \times \mu + C_{16} \times P^2 \times T + C_{17} \times \mu^3 + C_{18} \times T^3 + C_{19} \times P^3]$      (11)

291   where $C_0 = -207.739284$; $C_1 = -201.432367$; $C_2 = 1.1500875$; $C_3 = 0.678161$; $C_4 =$

292   $1.834310$; $C_5 = -1.309668$; $C_6 = -0.002251$; $C_7 = -25.879322$; $C_8 = -0.00201$; $C_9 =$

293   $-0.011747$; $C_{10} = 0.004118$; $C_{11} = 0.038415$; $C_{12} = -0.003692$; $C_{13} = 0.082021$; $C_{14} =$

294   $-2.664159 \times 10^{-7}$; $C_{15} = 0.001794$; $C_{16} = 3.978117 \times 10^{-5}$; $C_{17} = 3.600267$; $C_{18} =$

295   $1.477156 \times 10^{-6}$; and $C_{19} = -2.412235 \times 10^{-5}$

296       After rearrangement, the best resulting correlation based on GEP is expressed as

297   follows:

298      $D_c = 10^{-9} \times \left[\frac{A_1}{\mu} + \frac{A_2}{\sqrt{\mu}} + A_3 \times \sqrt{T} + A_4 \times P + A_5\right]$      (12)

299   The terms $A_1, A_2, A_3, A_4,$ and $A_5$ are defined as shown below:

$A_1 = -0.0001564 \times P^3 + 0.01113 \times P^2 + 0.02935 \times P - 2.83 \times \sqrt{P} + 8.362$

$A_2 = 0.02426 \times (P + T) - 2.466 \times \sqrt{P}$

$A_3 = 0.4583 + 0.123 \times P$

$$A_4 = 6.832 \times 10^{-5} \times P^2 + 0.003955 \times P + 19.34 \times \frac{1}{\sqrt{P}} - 3.801$$

$$A_5 = \frac{30.95 \times \ln(P)}{5.629 \times (P - \mu)} - \frac{0.0006024 \times T \times \sqrt{\mu}}{P} - 4.259 \times \ln(P^2 \times T) - 7.945$$

300    Fig. 6 illustrates the cross plots of the proposed correlations. The two cross plots

301  demonstrate a promising consistency between predictions of the correlations and the real data,

302  as very satisfactory alignments nearby the unit slope line are noticed for both GMDH and

303  GEP predictions. The fit is surprisingly trustworthy for the two correlations for both training

304  and testing data.

305    Furthermore, Fig. 7 shows the distribution of the percentage relative error between the

306  real values of the diffusivity coefficient of the $CO_2$ in brine and the values predicted by

307  GMDH and GEP correlations. According to this figure, the distributions of the relative error

308  around the zero-error line for both correlations at the training and testing phases are deemed

309  satisfactory. However, the reported results in this figure reveal that GEP correlation seems to

310  have a better predictive capability compared to GMDH correlation. For a detailed

311  quantification of the performances, Table 4 states the statistical criteria, namely AARD, $R^2$

312  and RMSE, for the proposed correlations. Moreover, a graphical comparison between the

313  global performances of the two correlations is reported in Fig. 8. It is clear from the

314  evaluation reported in Table 4 that the newly proposed correlations exhibit promising

315  predictive capabilities with overall values of AARD of 8.0404% and 4.3015% for GMDH and

316  GEP, respectively. Besides, the overall determination coefficient of the models indicates the

317  trustworthiness of the correlations fit. The analyses shown in Figs. 6–8 and Table 4 claim the

318  superiority of the GEP-based correlation in the prediction of the diffusivity coefficient of $CO_2$

319  in brine by considering the whole employed data points. Therefore, GEP correlation is used

320  for the rest of this paper.

321      The impact of the employed independent variables, namely pressure, temperature and

322   viscosity on the group error of the GEP-based correlation for prediction of the diffusivity

323   coefficient is investigated in Fig. 9 for the whole database. In Fig. 9(a), the GEP correlation

324   shows its worst predictions with an AARD of 5.31% where pressure is in the range of 20–40

325   MPa, while other intervals of pressure are predicted with an AARD value of less than 4.6%.

326   In Fig. 9(b), the maximum AARD value of the GEP correlation is obtained when temperature

327   is in the range of 298-323 K, while the rest of the intervals are estimated with an AARD value

328   that does not exceed 4.6%. In Fig. 9(c), the values of the diffusivity coefficient for viscosity

329   of less than 0.3 mPa.s are predicted with an AARD value of about 1%, while the AARD

330   values of other viscosities are between 4.5 and 5.5%. Consequently, the performance of the

331   proposed GEP is deemed again very sufficient in predicting the diffusivity coefficient.

332   **4.2. Comparison of the developed GEP correlation with decision trees (DTs), random**

333       **forest (RF) and the prior models**

334      It would be of interest from the reliability perspective to compare the results of the

335   proposed GEP correlation with other soft computing techniques, namely decision trees (DTs)

336   and random forest (RF) as well as the prior paradigms reported for the prediction of the

337   diffusivity coefficient of $CO_2$ in brine. To keep the work concise, more details about DTs and

338   RF can be found in published literature (Breiman, 2017; Guo et al., 2011; Peters et al., 2007;

339   Wilkinson, 2004). The final tuned control parameters of RF and DTs models are reported as

340   follows:

341       • RF: number of grown trees: 20; min leaf size: 5; min parent size: 2 × min leaf size;

342   predictor selection: interaction-curvature; splitting criterion: MSE.

343       • DTs: min leaf size: 1; min parent size: 5; quadratic error tolerance: 1E-6; predictor

344   selection: all-splits; prune criterion: MSE.

345    While developing DTs and RF models, 80% of the database was used for their training

346    and 20% of testing. The performance evaluation of the best-result DTs and RF paradigms

347    after various runs is reported in Table 5. By comparing the stated results in Tables 4 and 5, it

348    can be deduced that GEP based correlation outperforms both DTs and RF models.

349    The comparison includes the empirical models, namely those of Othmer and Thakar

350    (1953), Wilke and Chang (1955) and Cadogan et al. (2014a). In addition to the empirical

351    models, the implemented GEP correlation was compared with one of the most recent

352    intelligent paradigms proposed by Feng et al. (2019) based on hybrid genetic algorithm and

353    mixed Kernels-based support vector machine. It is worth mentioning that when performing

354    the comparison with the pre-existing approaches, we included only the points that satisfy the

355    applicability conditions in each correlation. Table 6 and Fig. 10 report the comparison of the

356    proposed GEP correlation and the models based on previously described statistical criteria.

357    According to these statistical analyses, among the existing models, the hybrid model proposed

358    by Feng et al. (2019) outperforms the available empirical models for predicting the diffusivity

359    coefficient with a global predictive AARD of 7.91%. Despite the exhibited accuracy of Feng

360    et al. (2019) model, it is worth mentioning that this model (based on mixed kernel SVR

361    coupled with GA) is resulted through performing some calculability efforts such as the

362    included quadratic programming involved in the establishment of the final solution. Besides,

363    as this paradigm is of the black box type, it is difficult to apply it to other related tasks.

364    The reported statistical quality measures in Table 6 and Fig. 10 demonstrate the

365    superiority of the newly proposed GEP correlation, as it outperforms both the prior intelligent

366    model and the empirical paradigms.

367    Being explicit based approaches, the performances of the empirical models and the

368    white-box GEP correlation were compared through the plot of the absolute relative error

369    distribution as shown in Fig. 11. As seen in this figure, 90% of the data points were predicted

15

370  by GEP correlation with an AARD value of less than 8.5%. The equivalent percentage of the

371  datapoints predicted with this AARD cutoff value by the empirical models are 50%, 25% and

372  24% for Cadogan et al. (2014a), Othmer and Thakar (1953), and Wilke and Chang (1955),

373  respectively.

374      As demonstrated in these comparative analyses, the newly implemented GEP

375  correlation was able to predict more reasonable values of the diffusivity coefficient of $CO_2$ in

376  brine. In addition, the GEP-based correlation has an explicit and simple form which can

377  predict the diffusivity coefficient of $CO_2$ in brine more directly than the other intelligent

378  schemes, and hence, it can be applied to other related tasks or implemented in different

379  softwares. The improvement brought by GEP based correlation in the prediction of the

380  coefficient of $CO_2$ diffusivity in brine can be explained by the followed learning strategy

381  during GEP steps which is based on the use of chromosomes, genes, functions, variables, and

382  the traditional and the new genetic operators, which result in more flexibility for capturing the

383  complexity of the modeled phenomenon.

384  **4.3. Trend Analysis**

385      To assess the efficiency of the implemented GEP correlation for accurate prediction of

386  diffusivity coefficient as a function of the input parameters, three different sets of

387  experimental measurements included in our study and the values predicted by GEP

388  correlation are depicted as function of the input parameters in Fig. 12. In Fig. 12(a), the

389  comparison is performed with respect to viscosity variation where pressure and temperature

390  are constant. In Fig. 12(b), the comparison is illustrated for different pressures and constant

391  temperature and viscosity. In Fig. 12(c), the comparison is shown for the case where pressure

392  is constant, and temperature and viscosity vary. Furthermore, additional comparison is

393  depicted in Fig. 13 by presenting the real measurements and the predictions of GEP as

16

function of pressure for the whole considered database. As shown in the plots of Figs. 12 and

13, the predicted and real diffusivity coefficient values of $CO_2$ in brine overlap properly

regardless of the operating conditions.

**4.4. Relative importance of input parameters**

A sensitivity analysis using the relevancy factor ($r$) (Chen et al., 2014; Hajirezaie et al.,

2015; Shateri et al., 2015), was performed to assess the relative importance of the input

variables on the diffusivity coefficient. The relevancy factor is expressed as follows:

$$r(I_j, O) = \frac{\sum_{i=1}^{n}(I_{j,i}-\overline{I_J})(o_i-\bar{o})}{\sqrt{\sum_{i=1}^{n}(I_{j,i}-\overline{I_J})^2 \sum_{i=1}^{n}(o_i-\bar{o})^2}} \qquad (13)$$

where the subscripts i and j refer to the data index and the variable, respectively; $I$ and $\bar{I}$

represent the input parameter and its average, respectively, while $O$ and $\bar{O}$ refer to the

predicted output and its average, respectively. It is worth noting that a high absolute ($r$) value

for an input parameter indicates its noteworthy impact on the output. Furthermore, achieving

positive/negative $r$ values for an input suggests a positive/negative effect on the output.

The obtained results regarding the relevancy factor for the diffusivity coefficient of $CO_2$

in brine are exhibited in Fig. 14. According to this figure, temperature has the biggest impact

on the outputs. In addition, it can be deduced that viscosity has a negative effect on the output,

while pressure and temperature positively affect the diffusivity coefficient.

**4.5. Outliers detection**

In the last part of this study, outliers detection was conducted to assess the quality of the

employed experimental data points employed for the establishment of the GEP correlation,

and also to define the applicability domain. The well-known Leverage approach was applied

(Rousseeuw and Leroy, 2005). The results from the Leverage approach are converted to the

famous graphical representation known as William plot (Nait Amar et al., 2019a, 2019b). This

plot scatters the standardized residual (R) of the predicted values versus the so-called hat (H)

418    values which corresponds to the diagonal elements of the hat matrix defined as (Gramatica,

419    2007; Rousseeuw and Leroy, 2005):

420    $$H = X(X^t X)^{-1} X^t \qquad (14)$$

421    where $X$ is a matrix with $(n \times d)$ size, with $n$ and $d$ represents the number of samples and

422    the variables, respectively, and $X^t$ is the transpose matrix of $X$. To delineate the applicability

423    in the Williams plot after presenting standardized residual as function of hat values, a

424    Leverage limit value (H*) calculated as $\frac{3(d+1)}{n}$ is utilized. In addition, the data points are

425    selected in the range of $\pm 3$ of standard deviation from the mean, where the cut-off value of 3

426    covers 99% of the distributed data (Gramatica, 2007; Rousseeuw and Leroy, 2005). The

427    suspected data points known as outliers are defined as the points which are situated in the

428    range of R > 3 or R < − 3 regardless of their hat value in comparison with H*. Hence,

429    existence of great accumulation of the data points in the ranges $0 \leq H \leq H^*$ and $-3 \leq R \leq 3$

430    indicates the high reliability of the model.

431        Fig. 15 shows the obtained Williams plot for the newly proposed correlation. This plot

432    reveals that 90 data points are in the intervals of $0 \leq H \leq 0.1304$ and $-3 \leq R \leq 3$, while

433    only two data points are found outside these margins, and hence, they are detected as doubtful

434    data. The Leverage approach confirms the statistical validity of the implemented GEP

435    correlation for predicting the diffusivity coefficient of $CO_2$ in brine.

## 5. Conclusions

437        In this paper, two new correlations were developed using GMDH and GEP for accurate

438    prediction of the diffusivity coefficient of $CO_2$ in brine. For developing the correlations, a

439    representative experimental database was collected from the published literature, based on

440    pressure, temperature and the viscosity of the solvent, as inputs. According to this study, the

441    following conclusions are drawn:

442    1. Both GMDH and GEP correlations showed very close prediction capabilities.

443    2. GEP correlation outperforms the GMDH correlation with an overall AARD value of

444       4.3014%.

445    3. The newly implemented GEP correlation exhibited very low AARD values with

446       respect to different intervals of input parameters.

447    4. The proposed GEP correlation can provide a fast and reasonably-priced estimation

448       of the coefficient of $CO_2$ diffusivity in brine.

449    5. The developed GEP correlation was compared with DTs, RF, mixed Kernels-based

450       support vector machine coupled with GA and other pre-existing models. The

451       accuracy of the developed correlation was superior to all these models.

452    6. The trends of the GEP outputs are logical in terms of the independent variables.

453    7. Temperature was found the most impacting parameter in the prediction of

454       diffusivity coefficient by GEP correlation.

455    8. The Leverage approach demonstrated the statistical validity of the model and only

456       two data points were detected as outliers.

**References**

459    Amini, S., Mohaghegh, S.D., Gaskari, R., Bromhal, G., others, 2012. Uncertainty analysis of a co2
460       sequestration project using surrogate reservoir modeling technique, in: SPE Western Regional
461       Meeting.
462    Aminu, K.T., McGlinchey, D., Chen, Y., 2019. Optimal design for real-time quantitative monitoring
463       of sand in gas flowline using computational intelligence assisted design framework. J. Pet. Sci.
464       Eng. 177, 1059–1071.
465    Azzolina, N.A., Nakles, D. V, Gorecki, C.D., Peck, W.D., Ayash, S.C., Melzer, L.S., Chatterjee, S.,
466       2015. CO2 storage associated with CO2 enhanced oil recovery: A statistical analysis of historical
467       operations. Int. J. Greenh. Gas Control 37, 384–397.
468    Bachu, S., Shaw, J.C., Pearson, R.M., others, 2004. Estimation of oil recovery and CO2 storage
469       capacity in CO2 EOR incorporating the effect of underlying aquifers, in: SPE/DOE Symposium
470       on Improved Oil Recovery.
471    Benamara, C., Nait Amar, M., Gharbi, K., Hamada, B., 2019. Modeling Wax Disappearance
472       Temperature Using Advanced Intelligent Frameworks. Energy & Fuels.
473       https://doi.org/10.1021/acs.energyfuels.9b03296
474    Bhakta, J.N., Lahiri, S., Pittman, J.K., Jana, B.B., 2015. Carbon dioxide sequestration in wastewater
475       by a consortium of elevated carbon dioxide-tolerant microalgae. J. CO2 Util. 10, 105–112.
476    Boot-Handford, M.E., Abanades, J.C., Anthony, E.J., Blunt, M.J., Brandani, S., Mac Dowell, N.,
477       Fernández, J.R., Ferrari, M.-C., Gross, R., Hallett, J.P., others, 2014. Carbon capture and storage

19

478    update. Energy Environ. Sci. 7, 130–189.
479  Breiman, L., 2017. Classification and regression trees. Routledge.
480  Cadogan, S., 2015. Diffusion of $CO_2$ in fluids relevant to carbon capture, utilisation and storage. PhD
481    thesis, Imp. Coll. London.
482  Cadogan, S.P., Hallett, J.P., Maitland, G.C., Trusler, J.P.M., 2015. Diffusion coefficients of carbon
483    dioxide in brines measured using 13C pulsed-field gradient nuclear magnetic resonance. J.
484    Chem. Eng. Data 60, 181–184. https://doi.org/10.1021/je5009203
485  Cadogan, Shane P, Hallett, J.P., Maitland, G.C., Trusler, J.P.M., 2014a. Diffusion coefficients of
486    carbon dioxide in brines measured using 13C pulsed-field gradient nuclear magnetic resonance.
487    J. Chem. Eng. Data 60, 181–184.
488  Cadogan, Shane P, Maitland, G.C., Trusler, J.P.M., 2014b. Diffusion coefficients of $CO_2$ and $N_2$ in
489    water at temperatures between 298.15 K and 423.15 K at pressures up to 45 MPa. J. Chem. Eng.
490    Data 59, 519–525.
491  Chen, G., Fu, K., Liang, Z., Sema, T., Li, C., Tontiwachwuthikul, P., Idem, R., 2014. The genetic
492    algorithm based back propagation neural network for MMP prediction in $CO_2$-EOR process.
493    Fuel 126, 202–212.
494  Choudhari, R. V., Doraiswamy, L.K., 1972. Physical Properties in Reaction of Ethylene and Hydrogen
495    Chloride in Liquid Media: Diffusivities and Solubilities. J. Chem. Eng. Data 17, 428–432.
496    https://doi.org/10.1021/je60055a012
497  Dargahi-Zarandi, A., Hemmati-Sarapardeh, A., Hajirezaie, S., Dabir, B., Atashrouz, S., 2017.
498    Modeling gas/vapor viscosity of hydrocarbon fluids using a hybrid GMDH-type neural network
499    system. J. Mol. Liq. 236, 162–171.
500  Davarazar, M., Jahanianfard, D., Sheikhnejad, Y., Nemati, B., Mostafaie, A., Zandi, S., Khalaj, M.,
501    Kamali, M., Aminabhavi, T.M., 2019. Underground carbon dioxide sequestration for climate
502    change mitigation--A scientometric study. J. $CO_2$ Util. 33, 179–188.
503  Ettehadtavakkol, A., Lake, L.W., Bryant, S.L., 2014. $CO_2$-EOR and storage design optimization. Int.
504    J. Greenh. Gas Control 25, 79–92.
505  Farajzadeh, R., Zitha, P.L.J., Bruining, J., 2009. Enhanced mass transfer of $CO_2$ into water:
506    experiment and modeling. Ind. Eng. Chem. Res. 48, 6423–6431.
507  Feng, Q., Cui, R., Wang, S., Zhang, J., Jiang, Z., 2019. Estimation of $CO_2$ diffusivity in brine by use
508    of the genetic algorithm and mixed kernels-based support vector machine model. J. Energy
509    Resour. Technol. 141, 41001.
510  Ferreira, C., 2001. Algorithm for solving gene expression programming: a new adaptive problems.
511    Complex Syst. 13, 87–129.
512  Frank, Marco J W, K.J. a M., Swaaij, W.P.M. Van, 1996. Marco J. W. Frank,* Johannes A. M.
513    Kuipers, and Wim P. M. van Swaaij. J. Chem. Eng. Data 41, 297–302.
514    https://doi.org/10.1021/je950157k
515  Frank, M.J.W., Kuipers, J.A.M., van Swaaij, W.P.M., 1996. Diffusion coefficients and viscosities of
516    $CO_2 + H_2O$, $CO_2 + CH_3OH$, $NH_3 + H_2O$, and $NH_3 + CH_3OH$ liquid mixtures. J. Chem. Eng.
517    Data 41, 297–302.
518  Gershenzon, N.I., Ritzi, R.W., Dominic, D.F., Soltanian, M., Mehnert, E., Okwen, R.T., 2015.
519    Influence of small-scale fluvial architecture on $CO_2$ trapping processes in deep brine reservoirs.
520    Water Resour. Res. 51, 8240–8256.
521  Gershenzon, N.I., Soltanian, M., Ritzi Jr, R.W., Dominic, D.F., 2014. Influence of small scale
522    heterogeneity on $CO_2$ trapping processes in deep saline aquifers. Energy Procedia 59, 166–173.
523  Gibbins, J., Chalmers, H., 2008. Carbon capture and storage. Energy Policy 36, 4317–4322.
524  Gozalpour, F., Ren, S.R., Tohidi, B., 2005. $CO_2$ EOR and storage in oil reservoir. Oil gas Sci.
525    Technol. 60, 537–546.
526  Gramatica, P., 2007. Principles of QSAR models validation: internal and external. QSAR Comb. Sci.
527    26, 694–701. https://doi.org/10.1002/qsar.200610151
528  Grude, S., Landrø, M., Dvorkin, J., 2014. Pressure effects caused by $CO_2$ injection in the Tubåen Fm.,
529    the Snøhvit field. Int. J. Greenh. Gas Control 27, 178–187.
530  Guo, L., Chehata, N., Mallet, C., Boukir, S., 2011. Relevance of airborne lidar and multispectral image
531    data for urban scene classification using Random Forests. ISPRS J. Photogramm. Remote Sens.
532    66, 56–66.

533  Guzmán, J., Garrido, L., 2012. Determination of carbon dioxide transport coefficients in liquids and
534      polymers by NMR spectroscopy. J. Phys. Chem. B 116, 6050–6058.
535  Hajirezaie, S., Hemmati-Sarapardeh, A., Mohammadi, A.H., Pournik, M., Kamari, A., 2015. A smooth
536      model for the estimation of gas/vapor viscosity of hydrocarbon fluids. J. Nat. Gas Sci. Eng. 26,
537      1452–1459.
538  Hemmati-Sarapardeh, A., Mohagheghian, E., 2017. Modeling interfacial tension and minimum
539      miscibility pressure in paraffin-nitrogen systems: Application to gas injection processes. Fuel
540      205, 80–89.
541  Hemmati-Sarapardeh, A., Varamesh, A., Husein, M.M., Karan, K., 2018. On the evaluation of the
542      viscosity of nanofluid systems: Modeling and data assessment. Renew. Sustain. Energy Rev. 81,
543      313–329.
544  Holtz, M.H., Nance, P.K., Finley, R.J., 2001. Reduction of greenhouse gas emissions through CO2
545      EOR in Texas. Environ. Geosci. 8, 187–199.
546  Ivakhnenko, A.G., 1971. Polynomial theory of complex systems. IEEE Trans. Syst. Man. Cybern.
547      364–378.
548  Jang, H.W., Yang, D., Li, H., 2018. A Power-Law Mixing Rule for Predicting Apparent Diffusion
549      Coefficients of Binary Gas Mixtures in Heavy Oil. J. Energy Resour. Technol. 140, 52904.
550  Jeong, H., Sun, A.Y., Lee, J., Min, B., 2018. A learning-based data-driven forecast approach for
551      predicting future reservoir performance. Adv. Water Resour. 118, 95–109.
552  Koza, J.R., 1992. Genetic programming II, automatic discovery of reusable subprograms. MIT Press,
553      Cambridge, MA.
554  Lee, J.W., Hawkins, B., Day, D.M., Reicosky, D.C., 2010. Sustainability: the capacity of smokeless
555      biomass pyrolysis for energy production, global carbon capture and sequestration. Energy
556      Environ. Sci. 3, 1695–1705.
557  Li, H., Yang, D., others, 2016. Determination of individual diffusion coefficients of solvent/CO 2
558      mixture in heavy oil with pressure-decay method. SPE J. 21, 131–143.
559  Lu, W., Guo, H., Chou, I.-M., Burruss, R.C., Li, L., 2013. Determination of diffusion coefficients of
560      carbon dioxide in water between 268 and 473 K in a high-pressure capillary optical cell with in
561      situ Raman spectroscopic measurements. Geochim. Cosmochim. Acta 115, 183–204.
562  Maharajh, D.M., 1975. Solubility and Diffusion of Gases in Water.
563  Maharajh, D.M., Walkley, J., 1972. The Temperature Dependence of the DiEusion Coefficients of Ar,
564      COz, CH4, CH3CI, CH3Br, and CHCB2F in Water. Can. J. Chem. 51, 944–952.
565  Mirjalili, S., 2015. How effective is the Grey Wolf optimizer in training multi-layer perceptrons. Appl.
566      Intell. 43, 150–161.
567  Mohagheghian, E., Hassanzadeh, H., Chen, Z., 2019. CO2 sequestration coupled with enhanced gas
568      recovery in shale gas reservoirs. J. CO2 Util. 34, 646–655.
569  Moultos, O.A., Tsimpanogiannis, I.N., Panagiotopoulos, A.Z., Economou, I.G., 2016. Self-diffusion
570      coefficients of the binary (H2O+ CO2) mixture at high temperatures and pressures. J. Chem.
571      Thermodyn. 93, 424–429.
572  Mutoru, J.W., Leahy-Dios, A., Firoozabadi, A., 2011. Modeling infinite dilution and Fickian diffusion
573      coefficients of carbon dioxide in water. AIChE J. 57, 1617–1627.
574  Nait Amar, M., Hemmati-Sarapardeh, A., Varamesh, A., Shamshirband, S., 2019a. Predicting
575      solubility of CO2 in brine by advanced machine learning systems: Application to carbon capture
576      and sequestration. J. CO2 Util. 33, 83–95.
577  Nait Amar, M., Zeraibi, N., 2019. An efficient methodology for multi-objective optimization of water
578      alternating CO2 EOR process. J. Taiwan Inst. Chem. Eng. 99, 154–165.
579  Nait Amar, M., Zeraibi, N., Hemmati-Sarapardeh, A., Shamshirband, S., 2019b. Modeling
580      temperature-based oil-water relative permeability by integrating advanced intelligent models
581      with grey wolf optimization: Application to thermal enhanced oil recovery processes. Fuel 242,
582      649–663.
583  Nait Amar, M., Zeraibi, N., Hemmati-Sarapardeh, A., Shamshirband, S., Mosavi, A., Chau, K., 2019c.
584      Modeling temperature dependency of oil-water relative permeability in thermal enhanced oil
585      recovery processes using group method of data handling and gene expression programming. Eng.
586      Appl. Comput. Fluid Mech. 13, 724–743.
587  Nijsing, R.A.T.O., Hendriksz, R.H., Kramers, H., 1959. Absorption of CO2 in jets and falling films of

588    electrolyte solutions, with and without chemical reaction. Chem. Eng. Sci. 10, 88–104.
589    https://doi.org/10.1016/0009-2509(59)80028-2

590  Nomeli, M.A., Riaz, A., 2017. A data driven model for the impact of IFT and density variations on
591    $CO_2$ storage capacity in geologic formations. Adv. Water Resour. 107, 83–92.

592  Othmer, D.F., Thakar, M.S., 1953. Correlating diffusion coefficient in liquids. Ind. Eng. Chem. 45,
593    589–593.

594  Peters, J., De Baets, B., Verhoest, N.E.C., Samson, R., Degroeve, S., De Becker, P., Huybrechts, W.,
595    2007. Random forests as a tool for ecohydrological distribution modelling. Ecol. Modell. 207,
596    304–318.

597  Piotrowski, A.P., Napiorkowski, J.J., 2012. Product-Units neural networks for catchment runoff
598    forecasting. Adv. Water Resour. 49, 97–113.

599  Reddy, K.A., Doraiswamy, L.K., 1967. Estimating liquid diffusivity. Ind. Eng. Chem. Fundam. 6, 77–
600    79. https://doi.org/10.1021/i160021a012

601  Riahi, K., Rubin, E.S., Taylor, M.R., Schrattenholzer, L., Hounshell, D., 2004. Technological learning
602    for carbon capture and sequestration technologies. Energy Econ. 26, 539–564.

603  Rostami, A., Hemmati-Sarapardeh, A., Karkevandi-Talkhooncheh, A., Husein, M.M., Shamshirband,
604    S., Rabczuk, T., 2019. Modeling heat capacity of ionic liquids using group method of data
605    handling: A hybrid and structure-based approach. Int. J. Heat Mass Transf. 129, 7–17.

606  Rousseeuw, P.J., Leroy, A.M., 2005. Robust regression and outlier detection. John wiley & sons.

607  Shahkarami, A., Mohaghegh, S., Gholami, V., Haghighat, A., Moreno, D., 2014. Modeling pressure
608    and saturation distribution in a $CO_2$ storage project using a Surrogate Reservoir Model (SRM).
609    Greenh. Gases Sci. Technol. 4, 289–315.

610  Shateri, M., Ghorbani, S., Hemmati-Sarapardeh, A., Mohammadi, A.H., 2015. Application of
611    Wilcoxon generalized radial basis function network for prediction of natural gas compressibility
612    factor. J. Taiwan Inst. Chem. Eng. 50, 131–141.

613  Tamimi, A., Rinker, E.B., Sandall, O.C., 1994. Diffusion Coefficients for Hydrogen Sulfide, Carbon
614    Dioxide, and Nitrous Oxide in Water over the Temperature Range 293–368 K. J. Chem. Eng.
615    Data 39, 330–332. https://doi.org/10.1021/je00014a031

616  Tan, K.K., Thorpe, R.B., 1992. Gas diffusion into viscous and non-Newtonian liquids. Chem. Eng.
617    Sci. 47, 3565–3572. https://doi.org/10.1016/0009-2509(92)85071-I

618  Teodorescu, L., Sherwood, D., 2008. High energy physics event selection with gene expression
619    programming. Comput. Phys. Commun. 178, 409–419.

620  Thomas, W.J., Adams, M.J., 1965. Measurement of the diffusion coefficients of carbon dioxide and
621    nitrous oxide in water and aqueous solutions of glycerol. Trans. Faraday Soc. 61, 668–673.
622    https://doi.org/10.1039/tf9656100668

623  Trevisan, L., Pini, R., Cihan, A., Birkholzer, J.T., Zhou, Q., Illangasekare, T.H., 2014. Experimental
624    investigation of supercritical $CO_2$ trapping mechanisms at the intermediate laboratory scale in
625    well-defined heterogeneous porous media. Energy Procedia 63, 5646–5653.

626  Venkatraman, V., Alsberg, B.K., 2017. Predicting $CO_2$ capture of ionic liquids using machine
627    learning. J. CO2 Util. 21, 162–168.

628  Versteeg, G.F., van Swaal, W.P.M., 1988. Solubility and Diffusivity of Acid Gases ($CO_2$, $N_2O$) in
629    Aqueous Alkanolamine Solutions. J. Chem. Eng. Data 33, 29–34.
630    https://doi.org/10.1021/je00051a011

631  Vivian, J.E., Peaceman, D.W., 1956. Liquid□side resistance in gas absorption. AIChE J. 2, 437–443.
632    https://doi.org/10.1002/aic.690020404

633  Wilke, C.R., Chang, P., 1955. Correlation of diffusion coefficients in dilute solutions. AIChE J. 1,
634    264–270.

635  Wilkinson, L., 2004. Classification and regression trees. Systat 11, 35–56.

636  Yan, Y., Xu, L., Lee, P., 2006. Mass flow measurement of fine particles in a pneumatic suspension
637    using electrostatic sensing and neural network techniques. IEEE Trans. Instrum. Meas. 55, 2330–
638    2334.

639  Yang, D., Tontiwachwuthikul, P., Gu, Y., 2006. Dynamic interfacial tension method for measuring
640    gas diffusion coefficient and interface mass transfer coefficient in a liquid. Ind. Eng. Chem. Res.
641    45, 4999–5008. https://doi.org/10.1021/ie060047e

642

643
644
645
646
647
648
649
650
651
652
653
654
655
656
657

| Solvent | Model | Expression | Included parameters |
|---------|-------|------------|---------------------|

658
659
660
661
662
663
664
665
666
667
668   **Table 1.** Summary of the existing empirical models for predicting the diffusivity coefficient of $CO_2$
669

23

| | | | |
|---|---|---|---|
| Brine | (Othmer and Thakar, 1953) | $$D_{CO_2} = \frac{14 \times 10^{-9}}{\mu^{1.1} V_m^{0.6}}$$ | • Molar volume of the diffusing substance ($V_m$ in cm³/gmol).<br>• Viscosity of the solvent ($\mu$ in $mPa \cdot s$) |
| | (Wilke and Chang, 1955) | $$D_{CO_2} = 7.4 \times 10^{-8} \frac{T\sqrt{\emptyset M}}{\mu V_m^{0.6}}$$ | • Temperature ($T$ in K).<br>• The association parameter $\phi$.<br>• Molecular weight of solvent (M). |
| | (Cadogan et al., 2014a) | $$D_{CO_2} = \frac{k_B T}{n_{SE} \pi \mu a}$$ | • $k_B$=1.38065×10⁻²³ J/K.<br>• $n_{SE}$ is the Stokes-Einstein number.<br>• The hydrodynamic radius of the solute (a in pm): $a = 168[1 + 2.0 \times 10^{-3}(T - 298)]$ |
| Pure water | (Lu et al., 2013) | $$D_{CO_2} = 13.942 \times 10^{-9} \left[\frac{T}{227} - 1\right]^{1.7094}$$ | • Temperature (T in K). |
| | (Moultos et al., 2016) | $$D_{CO_2} = D_0(P) \left[\frac{T}{Ts} - 1\right]^{m(P)}$$ | • $D_0$=$a_1$ln(P)+$a_2$, m=$b_1$ln(P)+$b_2$, where $a1$=-2.3097×10⁻⁹, $a2$=2,1064×10⁻⁸, $b1$=-0.17812 and $b2$=2.59406; $P$ is the pressure. |

670
671
672
673
674
675
676
677
678
679

680

681
682
683
684
685
686
687
688

24

689
690
691
692

**Table 2.** Summary of the gathered data

|  | Max | Avg. | Min | SD |
|---|---|---|---|---|
| **P (MPa)** | 49.30 | 9.64 | 0.1000 | 14.8030 |
| **T (K)** | 473.15 | 317.76 | 273 | 39.8 |
| **Viscosity (mPa.s)** | 1.9500 | 0.9003 | 0.1390 | 0.4720 |
| **Diffusivity coefficient ($\times 10^{-9}$ m$^2$/s)** | 16.1000 | 3.3522 | 0.3100 | 3.0874 |

693

694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734

25

735
736
737

**Table 3.** GEP setting parameters used in the study

| Parameters | Value/setting |
|---|---|
| Chromosome | 100 |
| Gene | 12 |
| Operators used | $+, -, \times, /, \exp., \sqrt{\phantom{x}}, INV, ln$ |
| Generations | 420 |
| Mutation rate | 0.45 |
| Inversion rate | 0.12 |

738

739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776

**Table 4.** Performance analysis of the implemented models.

| | | GEP | GMDH |
|---|---|---|---|
| **Training data** | AARD (%) | 3.8584 | 8.6269 |

26

| | | | |
|---|---|---|---|
| | R$^2$ | 0.9980 | 0.9943 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1427 | 0.2479 |
| **Test data** | AARD (%) | 6.0035 | 5.6292 |
| | R$^2$ | 0.9978 | 0.9874 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1245 | 0.2271 |
| **All data** | AARD (%) | 4.3014 | 8.0404 |
| | R$^2$ | 0.9979 | 0.9937 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1391 | 0.2440 |

**Table 5.** Performance analysis of the implemented decision trees (DTs) and random forest (RF) models.

| | | DTs | RF |
|---|---|---|---|
| **Training data** | AARD (%) | 4.2969 | 6.3627 |

27

| | | | |
|---|---|---|---|
| | $R^2$ | 0.9980 | 0.9973 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1598 | 0.1647 |
| **Test data** | AARD (%) | 8.8426 | 9.0015 |
| | $R^2$ | 0.9924 | 0.9940 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.2532 | 0.2764 |
| **All data** | AARD (%) | 5.1862 | 6.8790 |
| | $R^2$ | 0.9969 | 0.9966 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1785 | 0.1870 |

**Table 6.** Comparison of the performances with prior models

| | GEP | Feng et al. | Othmer | Wilke | Cadogan |
|---|---|---|---|---|---|

28

| | | | and Thakar | and Chang | et al. |
|---|---|---|---|---|---|
| **AARD (%)** | 4.3014 | 7.91 | 12.75 | 12.60 | 13.84 |
| **R$^2$** | 0.9979 | 0.9960 | 0.9661 | 0.9434 | 0.9858 |
| **RMSE** | 0.1391 | 0.1954 | 0.5661 | 0.7311 | 0.3661 |

903



904
905 **Fig. 1.** Frequency histograms of the collected dataset
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936

938          **Fig. 2.** Variation of diffusivity coefficient versus the independent variables

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954        **Fig. 3.** Results obtained for ten realizations with different GMDH highest orders

955

956

957

958

959

960

961

962

963

964

965

**Fig. 4.** Results of the GEP and GMDH sensitivity analysis on training and test sets

**Fig. 5.** A schematic structure of the implemented GMDH for predicting diffusivity coefficient

1016
1017
1018
1019
1020
1021



1022

**Fig. 6.** Cross plots of the established GEP and GMDH correlations for diffusivity coefficient prediction

1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055

35

1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066



1067
1068

1069

**Fig. 7.** Error distribution for the developed correlations

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

(a)



1087

(b)



1088

1089        **Fig. 8.** Comparison between the performances of the correlations: (a) AARD and (b) $R^2$

1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

37

(a)



1100

(b)



1101

(c)



1102

1103      **Fig. 9.** Comparison of the error distribution of the models with respect to input parameters

1104

1105

(a)



1106

(b)



1107

1108    **Fig. 10.** Comparison between the performances of GEP correlation and the prior models: (a)

1109    AARD and (b) $R^2$

1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120

39

1121
1122
1123

1124



1125

**Fig. 11.** Cumulative frequency vs. absolute percent relative deviation of GEP correlation and the prior empirical models.

1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154

40

**Fig. 12.** Comparison of the diffusivity coefficient obtained from measurements and generated by GEP correlation, as function of the input parameters.

1160



1161

1162  **Fig. 13.** Comparison of the diffusivity coefficient obtained from measurements and generated by GEP
1163                    correlation, for the whole employed pressure values.

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185



1186

1187 **Fig. 14.** Relevancy factor

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217

1218
1219
1220
1221
1222



1223

1224 **Fig. 15.** The Williams plot of GEP correlation.

1225
1226

1227

**Table 1.** Summary of the existing empirical models for predicting the diffusivity coefficient of $CO_2$

| Solvent | Model | Expression | Included parameters |
|---|---|---|---|
| Brine | (Othmer and Thakar, 1953) | $D_{CO_2} = \dfrac{14 \times 10^{-9}}{\mu^{1.1} V_m^{0.6}}$ | • Molar volume of the diffusing substance ($V_m$ in cm³/gmol). <br> • Viscosity of the solvent ($\mu$ in $mPa \cdot s$) |
| | (Wilke and Chang, 1955) | $D_{CO_2} = 7.4 \times 10^{-8} \dfrac{T\sqrt{\emptyset M}}{\mu V_m^{0.6}}$ | • Temperature ($T$ in K). <br> • The association parameter $\phi$. <br> • Molecular weight of solvent (M). |
| | (Cadogan et al., 2014a) | $D_{CO_2} = \dfrac{k_B T}{n_{SE} \pi \mu a}$ | • $k_B$=1.38065×10⁻²³ J/K. <br> • $n_{SE}$ is the Stokes-Einstein number. <br> • The hydrodynamic radius of the solute (a in pm): $a = 168[1 + 2.0 \times 10^{-3}(T - 298)]$ |
| Pure water | (Lu et al., 2013) | $D_{CO_2} = 13.942 \times 10^{-9} \left[\dfrac{T}{227} - 1\right]^{1.7094}$ | • Temperature (T in K). |
| | (Moultos et al., 2016) | $D_{CO_2} = D_0(P) \left[\dfrac{T}{Ts} - 1\right]^{m(P)}$ | • $D_0$=$a_1$ln(P)+$a_2$, m=$b_1$ln(P)+$b_2$, where $a1$=-2.3097×10⁻⁹, $a2$=2,1064×10⁻⁸, $b1$=-0.17812 and $b2$=2.59406; $P$ is the pressure. |

**Table 2.** Summary of the gathered data

|  | Max | Avg. | Min | SD |
|---|---|---|---|---|
| **P (MPa)** | 49.30 | 9.64 | 0.1000 | 14.8030 |
| **T (K)** | 473.15 | 317.76 | 273 | 39.8 |
| **Viscosity (mPa.s)** | 1.9500 | 0.9003 | 0.1390 | 0.4720 |
| **Diffusivity coefficient ($\times 10^{-9}$ m$^2$/s)** | 16.1000 | 3.3522 | 0.3100 | 3.0874 |

**Table 3.** GEP setting parameters used in the study

| Parameters | Value/setting |
|---|---|
| Chromosome | 100 |
| Gene | 12 |
| Operators used | $+, -, \times, /, \exp., \sqrt{\ \ }, INV, ln$ |
| Generations | 420 |
| Mutation rate | 0.45 |
| Inversion rate | 0.12 |

**Table 4.** Performance analysis of the implemented models.

| | | GEP | GMDH |
|---|---|---|---|
| **Training data** | AARD (%) | 3.8584 | 8.6269 |
| | $R^2$ | 0.9980 | 0.9943 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1427 | 0.2479 |
| **Test data** | AARD (%) | 6.0035 | 5.6292 |
| | $R^2$ | 0.9978 | 0.9874 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1245 | 0.2271 |
| **All data** | AARD (%) | 4.3014 | 8.0404 |
| | $R^2$ | 0.9979 | 0.9937 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1391 | 0.2440 |

**Table 5.** Performance analysis of the implemented decision trees (DTs) and random forest (RF) models.

| | | DTs | RF |
|---|---|---|---|
| **Training data** | AARD (%) | 4.2969 | 6.3627 |
| | $R^2$ | 0.9980 | 0.9973 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1598 | 0.1647 |
| **Test data** | AARD (%) | 8.8426 | 9.0015 |
| | $R^2$ | 0.9924 | 0.9940 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.2532 | 0.2764 |
| **All data** | AARD (%) | 5.1862 | 6.8790 |
| | $R^2$ | 0.9969 | 0.9966 |
| | RMSE ($\times 10^{-9}$ m$^2$/s) | 0.1785 | 0.1870 |

**Table 6.** Comparison of the performances with prior models

| | GEP | Feng et al. | Othmer and Thakar | Wilke and Chang | Cadogan et al. |
|---|---|---|---|---|---|
| **AARD (%)** | 4.3014 | 7.91 | 12.75 | 12.60 | 13.84 |
| **$R^2$** | 0.9979 | 0.9960 | 0.9661 | 0.9434 | 0.9858 |
| **RMSE** | 0.1391 | 0.1954 | 0.5661 | 0.7311 | 0.3661 |

**Fig. 1.** Frequency histograms of the collected dataset

**Fig. 2.** Variation of diffusivity coefficient versus the independent variables

**Fig. 3.** Results obtained for ten realizations with different GMDH highest orders

**Fig. 4.** Results of the GEP and GMDH sensitivity analysis on training and test sets

**Fig. 5.** A schematic structure of the implemented GMDH for predicting diffusivity coefficient

**Fig. 6.** Cross plots of the established GEP and GMDH correlations for diffusivity coefficient prediction

**Fig. 7.** Error distribution for the developed correlations

(a)



(b)



**Fig. 8.** Comparison between the performances of the correlations: (a) AARD and (b) $R^2$

**Fig. 9.** Comparison of the error distribution of the models with respect to input parameters

**Fig. 10.** Comparison between the performances of GEP correlation and the prior models: (a)

AARD and (b) $R^2$

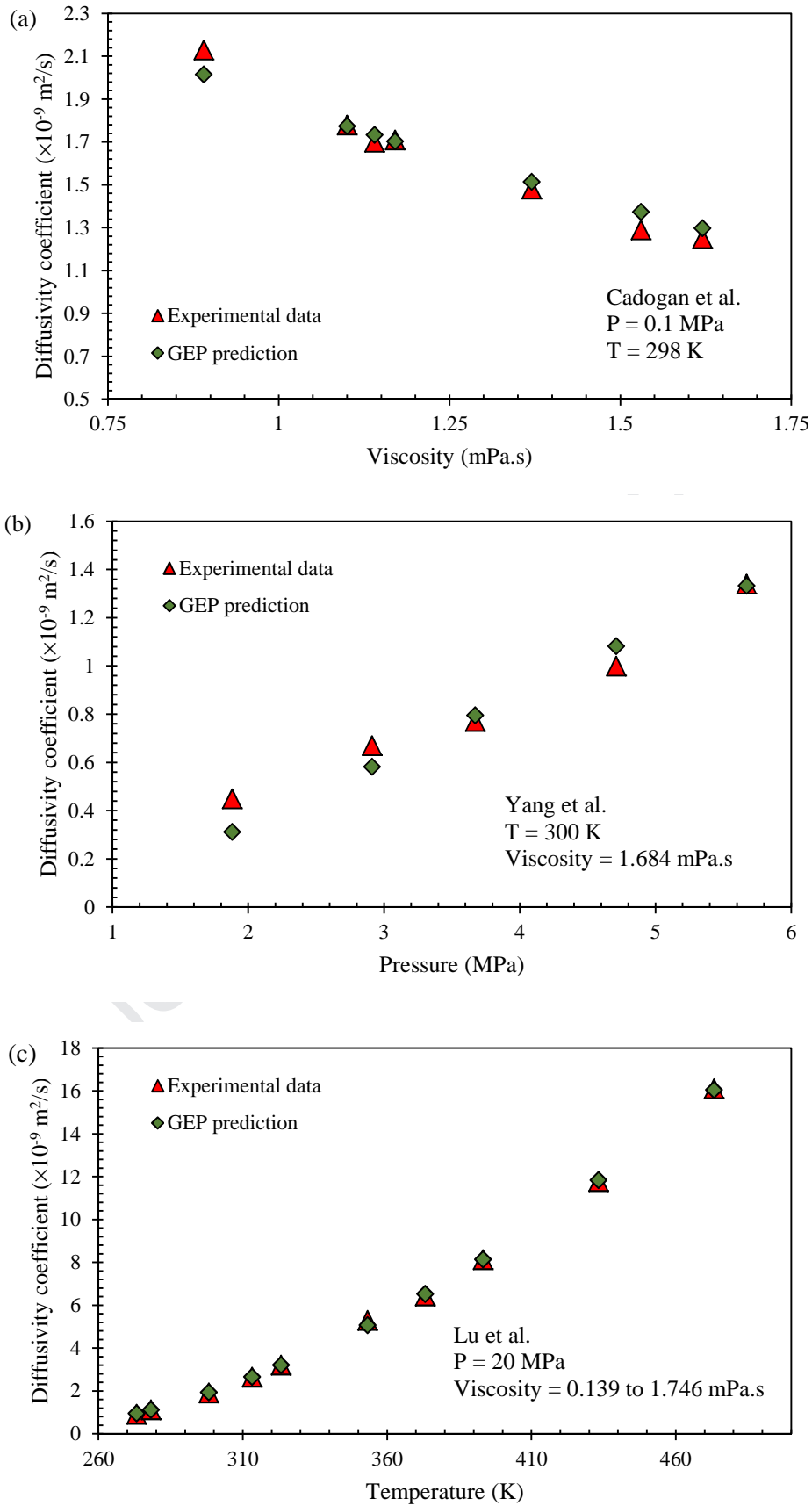**Fig. 11.** Cumulative frequency vs. absolute percent relative deviation of GEP correlation and the prior
empirical models.

**Fig. 12.** Comparison of the diffusivity coefficient obtained from measurements and generated by GEP correlation, as function of the input parameters.
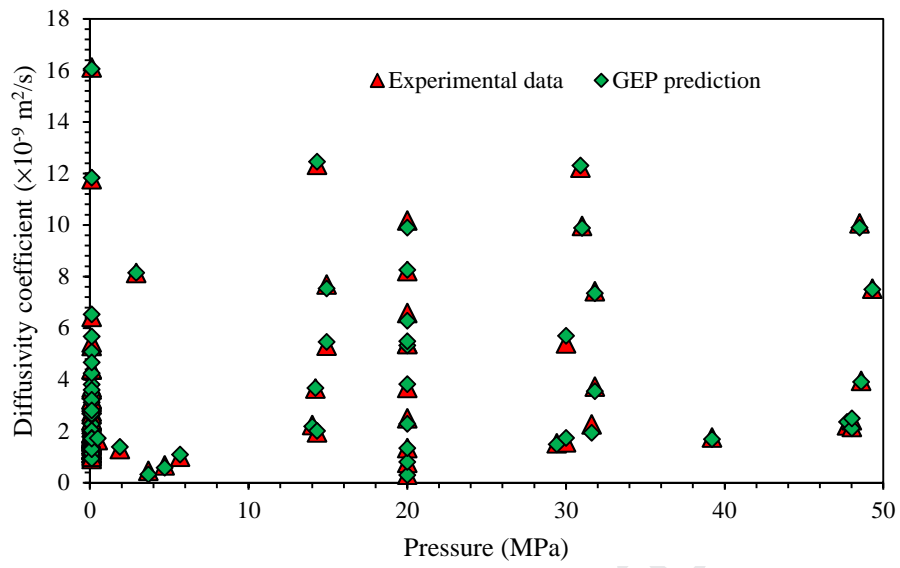
**Fig. 13.** Comparison of the diffusivity coefficient obtained from measurements and generated by GEP correlation, for the whole employed pressure values.
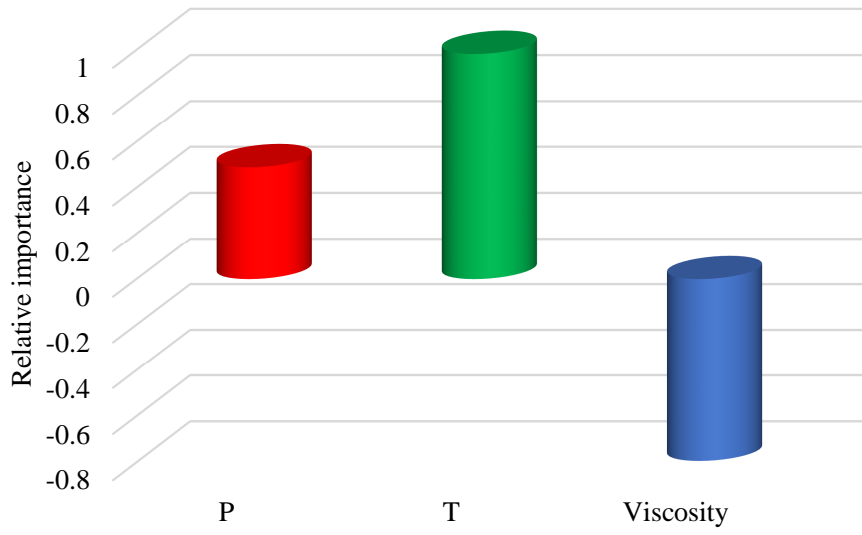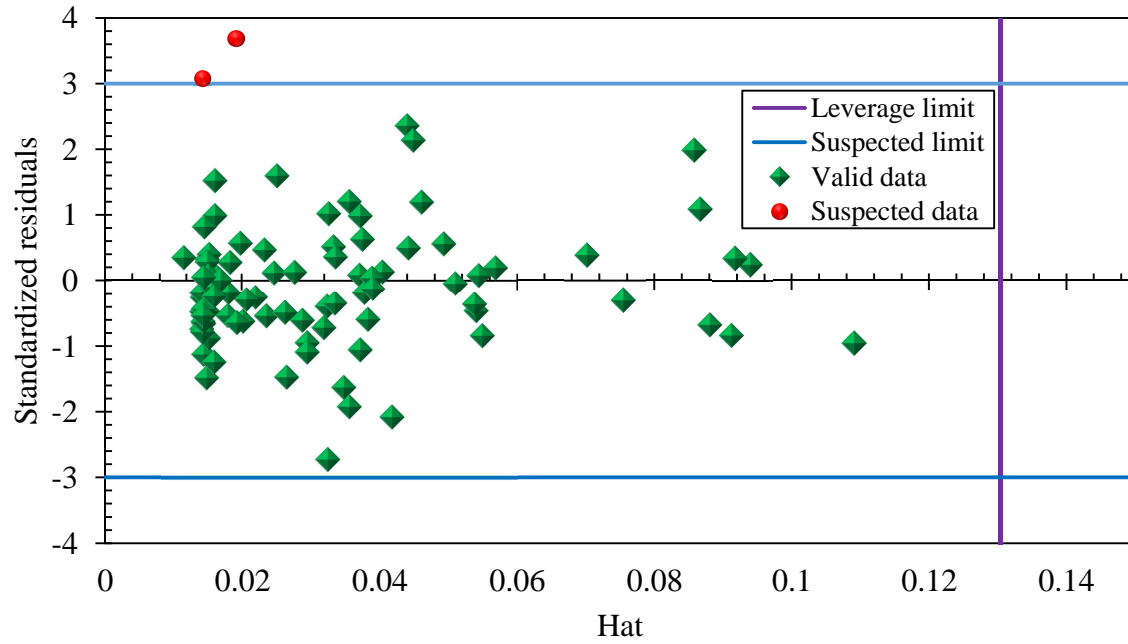
**Fig. 14.** Relevancy factor

**Fig. 15.** The Williams plot of GEP correlation.

# Highlights

- Two white-box machine learning techniques were implemented for predicting the diffusivity of $CO_2$ in brine.
- GEP is the best developed correlation.
- GEP correlation outperforms the prior paradigms.

# Authors' contributions

Menad Nait Amar: Data curation, Formal analysis, Methodology, Investigation and Modeling, Software, Writing.

Ashkan Jahanbani Ghahfarokhi: Supervision, Methodology, Writing, Reviewing and Editing.

**Declaration of interests**

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: