

Journal Pre-proof

Machine learning models for the prediction of diffusivities in supercritical CO₂ systems

José P.S. Aniceto, Bruno Zêzere, Carlos M. Silva



PII: S0167-7322(21)00007-6

DOI: <https://doi.org/10.1016/j.molliq.2021.115281>

Reference: MOLLIQ 115281

To appear in: *Journal of Molecular Liquids*

Received date: 30 November 2020

Revised date: 30 December 2020

Accepted date: 3 January 2021

Please cite this article as: J.P.S. Aniceto, B. Zêzere and C.M. Silva, Machine learning models for the prediction of diffusivities in supercritical CO₂ systems, *Journal of Molecular Liquids* (2021), <https://doi.org/10.1016/j.molliq.2021.115281>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2021 Published by Elsevier.

Machine learning models for the prediction of diffusivities in supercritical CO₂ systems

José P.S. Aniceto, Bruno Zêzere, and Carlos M. Silva

CICECO, Department of Chemistry, University of Aveiro, 3810-193 Aveiro, Portugal

Abstract

The molecular diffusion coefficient is fundamental to estimate dispersion coefficients, convective mass transfer coefficients, etc. Since experimental diffusion data is scarce, there is significant demand for accurate models capable of providing reliable diffusion coefficient estimations.

In this work we applied machine learning algorithms to develop predictive models to estimate diffusivities of solutes in supercritical carbon dioxide. A database of experimental data containing 13 properties for 174 binary systems totaling 4917 data points was used in the training of the models. Five machine learning algorithms were evaluated and the results were compared with three commonly used classic models.

The best results were found using the Gradient Boosted algorithm which showed an average absolute relative deviation (AARD) of 2.58 % (pure prediction). This model has five parameters: temperature, density, solute molar mass, solute critical pressure and solute acentric factor. For the same dataset, the classic Wilke-Chang equation showed AARD of 12.41 %.

1 Introduction

The knowledge of transport properties is required for the design, simulation and scale-up of rate controlled separations and chemical reactions. The binary diffusivities at infinite dilution, D_{12} , are fundamental to estimate important quantities like dispersion coefficients, convective mass transfer coefficients, and catalysts efficiency factors [1–3]. Over the past few years the so-called "green solvents" have been gaining more attention in both academia and industry, and from them we may detach supercritical carbon dioxide (SC-CO₂) [4, 5]. This solvent has been extensively used in

supercritical extraction (SFE), namely for extraction of compounds from vegetable matrices [5]. In this context, diffusivity coefficients data become extremely important. However, the knowledge of D_{12} in SC-CO₂ is still limited in terms of solutes and operating conditions, requiring accurate models capable of interpolating D_{12} and also predicting this property when no data is available [6].

Among the most well known models one can cite the hydrodynamic equation of Wilke-Chang [7], published in 1955, which is still widely used due to its simplicity since it only requires information on solvent viscosity, solute molecular mass, solute volume at normal boiling point and operating conditions. This equation has also been modified in several occasions giving rise, for instance, to the Lai-Tan equation [8], an empirical modification of the first one, specifically devised for SC-CO₂ systems. Furthermore, some other hydrodynamic based equations have been proposed and published for D_{12} estimation in SC-CO₂, which will not be addressed in this work but can be found elsewhere [9, 10]. Regarding correlative models, one can cite the 2-parameter correlation of Dymond-Hildebrand and Gatschinski (DHB) [11–13], based on free-volume theory, which is specially useful when some data of a given system is available allowing to both interpolate and extrapolate data for the desired condition.

Recently Artificial Intelligence has been applied in several fields of chemical engineering, for instance, for the estimation of physical properties of various compounds. Artificial neural network (ANN) models have been used to calculate the diffusion coefficients of pure compounds in water [14, 15]. Vaz *et al.* [16] proposed an ANN based correlation for the estimation of self-diffusion coefficients as a function of residual entropy and a molecular chain length parameter, which provided an average absolute relative deviation of 9.13 % for a large database with molecular dynamic and experimental values for hard-sphere, Lennard-Jones, hard-sphere chain, and real fluids composed of polar, nonpolar, symmetrical and asymmetrical molecules. Feed forward neural networks have also been developed to estimate the Fick diffusion coefficient in binary liquid systems [17]. Likewise, Eslamloueyan and Khademi [18] proposed a method based on a feed forward three-layer neural network to predict binary diffusion coefficient of gases at atmospheric pressure based on the critical temperature, critical volume and molecular weight of each component in the mixture. A genetic function approximation (GPA) derived model containing five parameters has been proposed to predict diffusion coefficient of non-electrolyte

organic compounds in air at 25 °C [19]. The estimation of the binary diffusion coefficients of liquid hydrocarbons at infinite dilution and in concentrated solutions has also been accomplished with Multi-layer perceptron (MLP) neural networks and an adaptive neuro-fuzzy inference system (ANFIS) [20].

In this work we applied machine learning algorithms, such as decision tree, nearest neighbors and ensemble methods to develop models for the prediction of binary diffusion coefficients of supercritical CO₂ systems. A large database of experimental data, covering small and large, polar and nonpolar solute molecules, was collected and used in the training of the models. The results were compared with the Wilke-Chang, Lai-Tan and DHB equations, extensively used for this purpose.

2 Theory, Database and Methods

The methodology used in this work to develop machine learning (ML) models for the prediction of diffusivities can be summarized in the following steps: (i) variable (feature) selection; (ii) learning algorithms selection; (iii) data splitting into training and testing sets; (iv) data scaling; (v) hyper-parameters optimization; and (vi) final model evaluation. These steps are detailed below. The ML models were compared with the classic models of Wilke-Chang [7], Lai-Tan [8] and Dymond-Hildebrand-Batschinski [11–13] shown in Section 2.4.

2.1 Database

The database used in this work has been updated and extended from the one previously published by Vaz *et al.* [10]. It is composed by 174 binary SC-CO₂ systems (solvent/solute) totaling 4917 data points and contains information on 13 properties, as shown in Table 1. It covers a wide range of temperatures and pressures: 283.15–398.15 K and 67–3500 bar, respectively.

Supercritical CO₂ densities were computed by the correlation of Pitzer and Schreiber [21] when they were not provided by the authors. The viscosities of SC-CO₂ were estimated by the correlation of Altunin and Sakhabetdinov [22] whenever necessary. The solute molar volumes at normal boiling point were estimated by Tyn–Calus equation [23] (Equation 4). The missing critical constants were estimated by Joback [23–25], Somayajulu [26], Klincewicz [27, 28], Ambrose [27, 29], and Constantinou-Gani [30] methods. The acentric factors, when not available, were estimated by the Lee-Kesler [31] and Pitzer [32] equations. The Lennard-Jones diameter and

energy were taken from Silva and Liu [13] and when not available were estimated by Equations 8 and 9 from Liu *et al.* [33].

Detailed information on the database used, including pure compound properties and reduced temperature and pressure ranges, is presented in the Supplementary Material Table SM1.

Table 1: Properties and variables available in the database of diffusivities of several solutes in SC-CO₂.

Property	Units	Description
D_{12}	$\text{cm}^2 \text{s}^{-1}$	Diffusion coefficient
T	K	Temperature
P	bar	Pressure
ρ_1	g cm^{-3}	Solvent density
μ_1	cP	Solvent viscosity
M_2	g mol^{-1}	Molar mass of solute
$T_{c,2}$	K	Critical temperature of solute
$T_{\text{bp},2}$	K	Boiling point temperature of solute
$P_{c,2}$	bar	Critical pressure of solute
$V_{c,2}$	mol cm^{-3}	Critical volume of solute
w_2		Acentric factor of solute
$\sigma_{\text{LJ},2}$	Å	Lennard-Jones diameter of solute
$\epsilon_{\text{LJ},2/k_B}$	K	Lennard-Jones energy of solute

2.2 Machine learning model development and optimization

Model features selection: Model features were selected from the properties and variables available in the database presented in Table 1. Features were excluded from the model until no collinearity above a defined threshold of 0.65 was present. To select the variable to exclude from a pair of collinear variables, the correlation with D_{12} and the ease of obtaining that variable were taken into account. Variables with low correlation with diffusivity were also removed.

Training and testing sets: In Machine Learning, data is usually divided into a training set, used for learning and fitting of the model, and a testing set, used to evaluate the fitted model after learning. Information from the testing set is never utilized during learning. Training and testing data sets were created by randomly splitting data base points 70 % into a training set and 30 % into testing sets. These data sets were kept unaltered for the evaluation of all models, guaranteeing the same input data for all.

Scaling: Feature scaling is usually beneficial to most learning algorithms as it often improves model robustness and training speed [34]. Scaling is accomplished through normalization or standardization of the features. Normalization consists in transforming the real range of values into a standard range (*e.g.* [0,1] or [-1,1]). Standardization consists in scaling variables so that they follow a standard normal distribution (mean of zero and standard deviation of one). In this work, properties were normalized to the [0,1] range using scikit-learn `MinMaxScaler`.

Hyper-parameter optimization: Unlike model parameters, which are fitted to the data when a model is trained, hyper-parameters are not learned from data and must be defined before training. They are configuration options of a given learning algorithm, usually with a numerical value, that influence how the model behaves. In this work, hyper-parameters were optimized by grid search with 4-fold cross-validation (using scikit-learn `GridSearchCV`). This method performs an exhaustive test of all hyper-parameters in a previously defined grid and evaluates the resulting model performance *via* k -fold cross-validation. The cross-validation technique avoids further reduction of the training set to create a validation set. Instead, the training set is split into k subsets and the model is trained using data from $k-1$ of the folds and tested on the remaining data. This is repeated using each $k-1$ combination of folds for training and the best hyper-parameters are those of the model with the best average performance. The tested hyper-parameters for each learning algorithm used, as well as the best hyper-parameters, are shown in Table 2.

Table 2: Tested and best hyper-parameter values for each machine learning algorithm. All remaining hyper-parameters were left at their default values.

ML Algorithm	Hyper-parameter	Values Tested	Best
--------------	-----------------	---------------	------

<i>k</i> -Nearest Neighbors	Number of neighbors	3; 4; 5; 6; 7; 10; 12; 15	3
	Algorithm	auto; ball_tree; kd_tree; brute	kd_tree
	Leaf size (BallTree or KDTree algorithm)	1; 2; 3; 5; 10; 15; 30	1
	Weight function	uniform; distance	distance
Decision Tree	Quality of a split metric	mse; mae	mae
	Split strategy	best; random	best
	Minimum number of samples per leaf	0.1; 1; 2; 5	1
	Minimum number of samples to split a node	0.1; 2; 4	2
	Minimum weighted fraction required for leaf node	0.1; 0.5	0
	Maximum number of features for split	auto; sqrt; log2; None	auto
	Minimum impurity decrease	0; 0.5; 2	0
Random Forest	Quality of a split metric	mse; mae	mae
	Number of estimators	5; 10; 15; 20; 30; 50; 100; 150	15
	Minimum number of samples per leaf	0.1; 1; 2; 5	1
	Minimum number of samples to split a node	0.1; 2; 4	4
	Minimum weighted fraction required for leaf node	0; 0.1; 0.5	0
	Maximum number of features for best split	auto; sqrt; log2; None	None

	Minimum impurity decrease	0; 0.5; 2	0
	Loss function	ls; lad; huber	lad
	Number of trees used in the boosting process	100; 500; 900; 1100; 1500	500
	Maximum depth of each tree	2; 3; 5; 10; 15	10
Gradient Boosted	Minimum number of samples per leaf	1; 2; 4; 6; 8	10
	Minimum number of samples to split a node	2; 4; 6; 10	1
	Maximum number of features for split	auto; sqrt; log2; None	sqrt

mse: mean squared error; mae: mean absolute error; ls: least square regression; lad: least absolute deviation; huber: a combination of ls and lad.

2.3 Machine Learning algorithms

Five ML algorithms were evaluated for the prediction of binary diffusivities: a Multilinear Regression, a k -Nearest Neighbors model, a Decision Tree algorithm, and two Ensemble Methods (Random Forest and Gradient Boosted). All models were implemented using the Python machine learning library scikit-learn version 0.22.1 [35]. In the following a brief description of each one is presented.

Multilinear Regression: A simple Ordinary Least Squares Multilinear Regression was used as a baseline model for the prediction of binary diffusivities. In a linear regression model, the real value, y , is a linear combination of features X_i weighted by coefficients b_i :

$$Y = b_0 + \sum_{i=1}^p X_i b_i \quad (1)$$

where Y is the predicted output, b_0 is the intercept or bias term, X_i are the input variables, b_i are the model coefficients, and p represents the number of parameters. The coefficients are optimized to minimize the residual sum of squares between the observed and the

calculated targets by the linear approximation [36]. This model was implemented using the *LinearRegression* class in scikit-learn.

Nearest Neighbors Regression: k -Nearest Neighbors (kNN) is a non-parametric method and one of the simplest machine learning algorithms. It operates by finding the k closest training examples (x_i) to each new input (x) and returns the average of their responses y_i .

$$Y(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i \quad (2)$$

where $N_k(x)$ is the neighborhood of x defined by the k closest points, where the Euclidean distance is usually used as a distance metric between samples [36]. Implementation was done using the *KNeighborsRegressor* in scikit-learn.

Decision Tree Regression: Decision Tree are models that use the training data to build a graph (tree) of simple decision rules that is used to analyse features. The prediction of the target variable is performed by following this tree, choosing the branches that return true values until an output (leaf) node is reached. This model was implemented using *DecisionTreeRegressor* in scikit-learn.

Ensemble methods: Occasionally, the algorithms presented above, due to their simplicity, cannot produce an accurate model for the problem at hand. In these cases other methodologies are required such as ensemble learning or neural networks. In ensemble methods a large number of simple models are trained and their predictions are combined to obtain a high accuracy model, providing improved generalizability and robustness over a single model [36]. There are two main types of ensemble methods. Averaging ensemble methods, like the Random Forest algorithm, average the predictions of several independently trained weak models. Boosting ensemble methods, such as the Gradient Boosted model, iteratively build multiple models in which each new learner mitigates the bias of the previous model. Random Forests and Gradient Boosted models are based on Decision Trees and have proven to be effective for regression on numerous cases [37]. Both ensemble methods were applied in scikit-learn using the *RandomForestRegressor* and *GradientBoostingRegressor* classes.

2.4 Classic D_{12} models

In this article the obtained ML models were compared with three widely known equations from the

literature, the predictive equations of Wilke-Chang [7] and Lai-Tan [8], and the correlation of Dymond-Hildebrand-Batschinski (DHB) [11–13]. The models are briefly described in the following.

Wilke-Chang equation: It is an empirical modification of the Stokes-Einstein relation and is given by:

$$D_{12} = 7.4 \times 10^{-8} \frac{(\phi M_1)^{0.5} T}{\mu_1 (V_{TC, bp, 2})^{0.6}} \quad (3)$$

where subscripts 1 and 2 represent solvent (CO₂) and solute, respectively, M_1 (g mol⁻¹) is the molar mass of solvent, μ_1 (cP) is the viscosity of the solvent, T (K) is the temperature, ϕ (nondimensional) is the association factor of the solvent (1.0 for the case of CO₂), and $V_{TC, bp, 2}$ (cm³ mol⁻¹) is the solute molar volume at normal boiling temperature. The last can be estimated by the Tyn-Calus relation [7, 38] when no experimental data is available:

$$V_{TC, bp, 2} = 0.285 \times V_{c, 2}^{1.48} \quad (4)$$

being $V_{c, 2}$ the critical volume of the solute

Lai-Tan equation: It is a modification of the Wilke-Chang equation and was specifically devised for tracer diffusion coefficients in SC-CO₂. It is described as:

$$D_{12} = 0.50 \times 10^{-7} \frac{M_1^{0.5} T}{(10 \times \mu_1)^{0.688} V_{c, 2}^{0.284}} \quad (5)$$

Dymond-Hildebrand Batschinski model (DHB): It is a free-volume based model frequently adopted to describe transport properties in nonpolar systems. The equation is [11–13]:

$$D_{12} = B_{DHB} \sqrt{T} (V_1 - V_D) \quad (6)$$

where V_1 (cm³ mol⁻¹) is the molar volume of the solvent, B_{DHB} (cm³ mol⁻¹ K^{-1/2}) is a parameter characteristic of the solute-solvent pair, and V_D (cm³ mol⁻¹) is the minimum solvent molar volume required for diffusion. The last two are the adjustable parameters of the equation.

2.5 Model evaluation

The performance of the proposed models was evaluated using the average absolute relative deviation (AARD), which was always calculated in order to assess the goodness of fittings and predictions. It is given by:

$$\text{AARD}(\%) = \frac{100}{\text{NDP}} \sum_{i=1}^{\text{NDP}} \left| \frac{D_{12}^{\text{calc}} - D_{12}^{\text{exp}}}{D_{12}^{\text{exp}}} \right|_i \quad (7)$$

where superscripts *calc* and *exp* denote calculated and experimental values, and NDP is the number of data points. In addition to this weighted average, the simple arithmetic average of the AARD values of all systems (AARD_{arithmetic}) was also calculated. The minimum and maximum system AARD are reported as an indication of the performance of the best and worst systems. Likewise, the AARD metric was also applied to the classic models used for comparison.

3 Results and discussion

Model development started with the selection of the relevant properties and variables from Table 1 (feature selection), followed by the selection of the machine learning algorithm and finally the comparison of the best machine learning models with the Wilke-Chang, Lai-Tan and DHB equations.

3.1 Properties and variables selection

A feature selection process was conducted to identify appropriate variables and properties for the model. Figure 1 shows the correlation matrix (in the form of heat map) for the SC-CO₂ data set, where values represent absolute Pearson correlation. Collinear quantities were excluded from the model by analyzing the Pearson correlations and setting a correlation coefficient threshold of 0.65. For each pair of quantities with a correlation above this value, usually the one with lower correlation with D_{12} was removed from the model. The simplicity of the model was also taken into account when selecting/excluding variables. For instance, when analyzing $T_{\text{bp},2}$ and M_2 pair, which presents a correlation of 0.88, $T_{\text{bp},2}$ has slightly higher correlation with D_{12} ; however, $T_{\text{bp},2}$ was excluded in favor of M_2 as information on M_2 is immediate and rigorously calculated, thus endowing the model with greater simplicity.

Figure 1: Pearson correlation heat map for all properties and variables available for the supercritical CO₂ model.

Following this procedure, the selected properties/variables for the model were temperature, pressure, density, solute molar mass, solute critical pressure and solute acentric factor. Upon further testing, pressure was also excluded as its effect upon the model performance was negligible. This is consistent with the fact that pressure and temperature effects are both included in the density effect, thus at least one of these variables should be theoretically unnecessary. Table 3 shows the final properties/variables chosen for the supercritical CO₂ machine learning model (ML SC-CO₂), as well as those embodied in the Wilke-Chang and Lai-Tan equations used for comparison. The DHB correlation relies on V_1 and two fitted parameters (B_{DHB} and V_D), presented in Section 2.4.

The widespread use of the classic models such as the Wilke-Chang equation owes to their ease of applicability. Wilke-Chang and Lai-Tan equations require variables and parameters which are simple to obtain as temperature, viscosity, critical volume and molecular mass. The ML model here proposed relies on similarly simple quantities (temperature, density, molecular mass, and critical pressure) but also on the acentric factor whose information is less easily available. However, this property can be estimated using either the Pitzer [32] or Lee-Kesler [31] equations. The former requires knowledge of $T_{c,2}$, $P_{c,2}$ and vapor pressure while the later requires information on $T_{c,2}$, $P_{c,2}$ and $T_{\text{bp},2}$.

Table 3: System properties/variables used in each model.

Parameters	Proposed model		Classic models	
	ML SC-CO ₂	Wilke-Chang	Lai-Tan	DHB
T	✓	✓	✓	
P				
ρ_1	✓			
μ_1		✓	✓	
M_2	✓			
$T_{c,2}$				
$T_{\text{bp},2}$				
$P_{c,2}$	✓			

$V_{c,2}$		✓	✓	
w_2	✓			
$\sigma_{LJ,2}$				
$\varepsilon_{LJ,2/k_B}$				
M_{CO_2}		✓	✓	
V_1				✓
B_{DHB}				✓ ^a
V_D				✓ ^a
Count	5	4	4	3

^a Requires fitting to experimental data.

3.2 Machine Learning model selection

Five machine learning algorithms were applied in this study covering several types of supervised learning models: Multilinear regression (linear model), k -Nearest Neighbors, Decision Tree, Random Forest (averaging ensemble method) and Gradient Boosted (boosting ensemble method).

In Figure 2 the predicted diffusivities are plotted *versus* the experimental values for the test set for the five machine learning algorithms. The best results are achieved using the Gradient Boosted algorithm (Figure 2e) with AARD = 2.58 %. Note this is the pure prediction deviation calculated for the test set, not used in model training. Remaining algorithms ranked, from lower to higher AARD, as: Random Forest (4.14 %), k -Nearest Neighbors (4.77 %), Decision Tree (4.89 %), and Multilinear Regression (15.81 %). Similar behavior is observed for the k -Nearest Neighbors, Decision Tree, Random Forest and Gradient Boosted models (Figures 2b, 2c, 2d, and 2e) which present random distribution along diagonal, while the Multilinear Regression shows heavy underfitting for low and high D_{12} values.

Each ML algorithm hyper-parameters were optimized as described in Section 2.2. A grid search with 4-fold cross-validation was applied, using at least 3 levels for numeric hyper-parameters and all available options for non-numeric hyper-parameters. A detailed description of the tested and best hyper-parameter values for each machine learning algorithm is provided in Table 2. In the case of the best model (Gradient Boosted), it uses the least absolute

deviation as the loss function; 500 boosting stages; a minimum number of samples required to split of 10; maximum depth of 10; leafs with a minimum of 1 sample; and a maximum number of features considered when splitting equal to the square root of the total number of features; all remaining hyper-parameters were left at their default values. The best model will be henceforth denoted by ML-GB SC-CO₂ (Machine Learning Gradient Boosted model applied to D_{12} in SC-CO₂).

Figure 2: Predicted *versus* experimental diffusivities for the test set using different machine learning algorithms: a) Multilinear Regression; b) k -Nearest Neighbors; c) Decision Tree; d) Random Forest; and e) Gradient Boosted.

3.3 Comparison with classic models

The final proposed model (ML-GB SC-CO₂) contains five parameters: temperature, solvent density, solute molar mass, solute critical pressure and solute acentric factor. It provided an AARD of 2.58 % for the test set, which contains 168 systems and 1476 data points.

Table 4 compares this model performance with the classic models of Wilke-Chang, Lai-Tan, and DHB, in terms of global AARD, arithmetic average of systems (AARD_{arithmetic}), as well as minimum and maximum system AARD (AARD_{min} and AARD_{max}). The DHB equation takes one less system since not enough points were available in the training set to fit its two parameters. Overall the new ML-GB SC-CO₂ model outperforms the classic models, with only the DHB equation attaining comparable results, which may be attributed to the two embodied parameters. With regard to the maximum system AARD, ML-GB SC-CO₂ shows a maximum deviation of 17.27 %, slightly below the one provided by DHB (19.54 %) and significantly below those of the remaining classic models. Additionally, one should keep in mind that the ML model is universal while the DHB equation is system-specific, *i.e.* requires two parameters previously fitted to available data of each system.

The performance of the classic models is presented in Figure 3 in terms of predicted *versus* experimental diffusivities. It is interesting to note that the Lai-Tan equation, a modification of the Wilke-Chang equation specifically developed for supercritical CO₂ systems, presents the worst results with an AARD of 26.01 %. This inferior performance is because the Lai-Tan equation was

obtained by re-optimizing the Wilke-Chang frontal coefficient and μ_1 and $V_{c,2}$ exponents using only 141 experimental points from 8 systems [8]. Consequently, such weaker support gives rise to higher errors when Lai-Tan equation deviates from its original pressure and temperature conditions and type of molecules.

Table 4: Comparison of the performance of the ML-GB SC-CO₂ model for the prediction (test set) of diffusivities in SC-CO₂ with the classic models: number of test systems, number of test points (NDP), global AARD, arithmetic average of systems AARD (AARD_{arithmetic}), minimum and maximum AARD (AARD_{min} and AARD_{max}).

Model	Systems	NDP	AARD (%)	AARD _{arithmetic} c (%)	AARD _{min} (%)	AARD _{max} (%)
ML-GB	168	1476	2.58	2.77	0.31	17.27
SC-CO ₂						
Wilke-Chan	168	1476	12.41	14.00	2.40	54.04
g						
Lai-Tan	168	1476	26.01	22.97	1.56	84.25
DHB	167	1473	4.27	4.03	0.26	19.54

Figure 3: Predicted *versus* experimental diffusivities for the test set using the (a) Wilke-Chang, (b) Lai-Tan, and (c) DHB equations.

Detailed results for each system in the test and train sets are presented in Table 5, specifying the solute, number of data points and global AARD, for the ML-GB SC-CO₂ model and the three classic models adopted for comparison. The best results were obtained for the n-butylbenzene (0.31 %), n-decane (0.34 %) and 1-naphthol (0.44 %) systems, while the worst ones correspond to the m-xylene system (17.27 %) followed by the 1-methylnaphthalene (10.53 %) and 1-hexadecene (9.02 %) systems.

Figures 4 and 5 illustrate the $D_{12} / T^{0.5}$ dependence on solvent molar volume (free-volume theory) and Stokes-Einstein plots, respectively, using experimental and predicted (ML-GB SC-CO₂ model) diffusivities for two systems: acetone and 1-methylnaphthalene. The main idea

behind such representations is to demonstrate that the ML-GB SC-CO₂ model conserves the classical trends verified in the Stokes-Einstein and free-volume representations. As can be observed in both figures, the expected trends are kept, conserving the linear dependency of $D_{12}/T^{0.5}$ versus V_1 and D_{12} versus $T\mu_1^{-1}$. In fact this dependency can be quantified by the coefficient of determination (R^2) obtained scoring a value of 0.9509 (acetone) and 0.9681 (1-methylnaphthalene) for the Stokes-Einstein relation, and 0.9772 (acetone) and 0.9443 (1-methylnaphthalene) for the free-volume plot. In the case of 1-methylnaphthalene one sees that the model still follows the expected linear trends in the Stokes-Einstein and free-volume plots, notwithstanding the larger deviations of some data, mainly the last point. This causes the higher AARD found for this system.

Table 5: Calculated results (AARD) for the diffusivities of solutes in supercritical CO₂ for every system in the test and train sets achieved by the ML-GB SC-CO₂ model and the classic models used for comparison. Systems sources and ranges of temperature, pressure and densities are reported in Table SM2 in Supplementary Material.

Solute	ML-GB			AARD (%)							
	SC-CO ₂			Wilke-Chan		Lai-Tan		DHB			
	Test	Train	g	Test	Train	Test	Train	Test	Train		
α -linolenic acid	56	13	43	1.39	0.69	13.19	14.56	30.9	32.7	2.62	2.83
								4	0		
α -pinene	30	13	17	3.17	1.54	6.51	5.91	13.7	13.4	3.10	2.91
								7	0		
α -tocopherol	82	28	54	1.47	0.77	25.86	26.75	30.8	31.6	1.96	2.23
								5	3		
β -carotene	90	25	65	1.35	0.79	15.79	14.53	66.9	65.8	2.36	2.30
								7	9		
β -pinene	15	2	13	4.87	1.35	5.03	12.75	4.65	5.14	10.3	3.62

γ -linolenic acid	142	37	105	1.23	0.81	7.63	7.84	36.3	36.3	2.40	2.08		
										5	6		
γ -linolenic acid ethyl ester	41	13	28	2.70	1.43	8.57	6.16	25.5	23.3	5.45	4.98		
										0	0		
γ -linolenic acid methyl ester	52	16	36	2.41	1.05	11.91	14.07	13.2	22.8	7.04	7.71		
										9	9		
1,1,1,5,5,5-hexafluoroacetylacetone	15	5	10	3.67	1.35	20.18	18.34	35.0	30.8	4.87	4.24		
										6	4		
1,1'-dimethylferrocene	68	25	43	1.71	0.77	11.39	12.32	18.7	18.6	3.91	3.73		
										2	9		
1,2-dichlorobenzene	15	4	11	2.60	0.80	7.27	6.75	19.8	16.7	3.05	1.91		
										2	2		
1,2-diethylbenzene	15	4	11	3.03	0.61	8.49	5.36	17.8	17.0	3.71	2.30		
										0	9		
1,3,5-trimethylbenzene	31	8	26	3.05	1.38	6.31	6.26	14.1	15.9	6.04	3.49		
										4	5		
1,3-dichlorobenzene	4	1	3	2.11	1.08	11.78	11.43	25.8	23.2	0.95	1.17		
										3	2		
1,3-divinylbenzene	15	6	9	2.10	0.74	5.14	2.80	18.5	15.3	1.60	1.52		
										1	0		
1,4-diethylbenzene	15	4	11	1.92	0.55	7.23	5.08	17.1	18.3	6.30	1.96		
										4	9		
15-crown-5	29	13	16	2.26	0.93	6.91	8.62	13.5	16.3	6.44	5.95		
										1	7		
1-hexadecene	11	6	5	9.02	0.75	16.52	15.75	16.9	26.1	10.3	14.8		
										1	9	7	1
1-methylnaphthalene	11	2	9	10.5	0.79	30.84	18.02	39.6	27.3	19.5	3.21		
						3				5	2	3	
1-naphthol	11	3	8	0.44	0.45	5.43	5.90	3.57	4.99	1.29	0.97		
1-phenyldodecane	15	5	10	1.16	0.90	6.15	7.37	51.2	45.1	2.31	3.98		
										4	4		

1-phenylethanol	15	2	13	3.45	0.46	16.31	9.39	25.4	25.3	7.94	1.99		
									4	1			
1-phenylhexane	15	6	9	1.21	0.64	8.23	7.05	23.1	24.1	2.79	2.88		
									6	9			
1-phenyloctane	15	1	14	1.13	0.78	11.44	8.45	20.1	27.8	6.10	3.37		
									7	6			
1-propanol	17	4	13	3.70	0.40	8.70	17.50	7.04	3.70	4.19	2.45		
2,2,4,4-tetramethyl-3-pentanone	9	2	7	6.29	0.94	29.08	26.41	26.4	19.1	1.38	0.64		
									6	8			
2,3-dimethylaniline	15	5	10	2.43	0.70	11.17	18.47	33.1	33.5	3.12	2.37		
									0	8			
2,4-dimethyl-3-pentanone	8	2	6	6.27	1.13	12.32	10.90	15.6	20.9	2.31	2.51		
									3	8			
2,4-dimethylphenol	15	6	9	1.83	0.53	5.66	11.92	22.4	25.9	4.85	2.91		
									3	6			
2,6-dimethylaniline	15	3	12	0.81	0.65	12.43	11.23	28.2	28.0	1.42	3.87		
									1	2			
2,6-dimethylnaphthalene	6	0	6	n.d.	3.36	n.d.	7.15	n.d.	18.6	n.d.	4.24		
									8				
2,7-dimethylnaphthalene	6	1	5	6.64	4.02	8.09	6.67	8.66	17.5	7.99	3.74		
									1				
2-bromoanisole	15	2	13	0.93	0.48	17.40	16.39	30.0	30.7	3.58	3.68		
									7	9			
2-butanone	40	13	27	1.70	0.47	3.60	5.87	5.24	4.46	2.34	3.08		
2-ethyltoluene	15	3	12	3.68	0.64	8.69	9.01	10.2	10.6	4.53	3.71		
									6	0			
2-fluoroanisole	15	6	9	2.47	0.48	21.21	16.65	25.3	27.5	3.59	1.81		
									7	1			
2-heptanone	11	1	10	6.24	0.52	34.47	29.67	32.8	22.1	1.05	1.93		
									4	0			
2-methylanisole	15	5	10	5.03	3.22	9.09	9.97	22.5	23.0	4.30	2.56		

								0	8		
2-naphthol	16	3	13	0.94	0.24	7.32	7.96	10.8	10.6	2.45	1.68
								1	9		
2-nitroanisole	15	3	12	2.54	0.74	12.55	11.20	33.2	29.9	4.15	2.25
								7	0		
2-nonanone	10	3	7	1.96	1.55	36.44	35.47	26.4	24.6	3.23	2.13
								7	6		
2-pentanone	23	4	19	4.09	0.48	7.92	3.71	3.76	1.96	5.31	1.71
2-phenyl-1-propanol	15	6	9	2.13	0.78	10.53	9.19	28.7	30.9	3.96	2.20
								3	3		
2-phenylethanol	15	4	11	1.02	0.22	4.16	11.34	24.0	27.7	3.89	2.68
								9	7		
2-phenylethyl acetate	15	10	5	3.07	0.42	8.19	9.52	36.1	36.7	4.34	1.80
								8	5		
2-propanol	15	7	11	1.31	0.52	10.50	8.98	6.15	8.59	3.71	1.65
3-ethyltoluene	15	9	6	2.08	0.53	10.60	13.31	8.85	9.46	7.21	3.12
3-fluorophenol	4	0	4	n.d.	0.61	n.d.	13.15	n.d.	24.2	n.d.	1.00
										6	
3-methylbutylbenzene	15	4	11	0.70	0.73	4.16	6.85	18.0	21.4	4.06	2.69
								8	6		
3-nitrotoluene	15	3	12	1.97	0.97	2.40	4.40	13.1	19.7	6.40	3.91
								0	4		
3-pentanone	46	12	34	1.59	0.57	7.76	8.47	4.95	3.97	2.78	2.59
3-phenyl-1-propanol	15	7	8	2.16	0.61	4.78	7.45	28.8	25.7	4.26	1.92
								0	9		
3-phenylpropyl acetate	15	5	10	1.75	0.72	3.92	8.72	39.5	37.3	5.21	3.37
								0	3		
4-ethyltoluene	15	2	13	2.71	0.41	10.97	6.92	10.4	13.0	4.90	2.78
								0	8		
4-heptanone	9	1	8	0.74	0.41	36.46	36.54	28.5	29.8	0.39	0.48
								8	2		

4-methylanisole	15	3	12	3.82	2.96	17.29	17.58	29.0	32.5	2.11	3.64		
										3	1		
5-nonanone	12	1	11	0.48	0.76	31.24	34.59	11.6	21.6	0.73	1.17		
										3	4		
5-tert-butyl-m-xylene	31	13	18	1.34	0.49	8.85	8.23	19.7	19.1	3.69	3.97		
										3	5		
6-undecanone	13	5	8	6.02	1.27	36.61	36.47	20.9	18.7	3.40	2.00		
										2	4		
AA ethyl ester	48	13	35	0.72	0.40	15.20	15.15	31.9	29.8	1.60	0.98		
										0	8		
acetone	213	74	139	2.91	1.18	4.89	5.99	10.2	10.5	5.44	5.38		
										0	4		
acridine	6	2	4	4.55	0.77	6.93	3.93	29.1	26.2	3.85	2.93		
										4	7		
adamantanone	5	2	6	3.02	0.45	12.37	20.10	13.5	15.4	11.0	0.89		
										3	3	0	
allylbenzene	15	8	7	3.32	0.44	6.62	3.93	17.8	14.6	5.53	1.76		
										0	4		
aniline	15	6	9	6.59	0.42	37.92	30.28	34.2	33.2	3.83	2.21		
										5	4		
anisole	15	3	12	1.48	0.59	6.25	7.60	16.4	16.2	1.90	3.28		
										7	7		
anthracene	22	8	14	3.48	0.67	10.61	10.25	12.9	15.6	1.87	1.31		
										4	6		
arachidonic acid	75	23	52	1.58	0.92	9.51	9.78	41.3	41.0	2.88	2.18		
										5	0		
behenic acid ethyl ester	17	5	12	0.64	0.60	21.03	21.47	33.1	30.5	0.93	0.84		
										0	6		
benzene	249	84	165	5.64	1.68	9.40	8.44	8.93	9.54	7.12	8.13		
benzoic acid	35	7	28	4.78	0.50	11.02	9.97	15.0	15.7	5.37	6.83		
										5	0		

benzyl acetate	15	6	9	1.29	0.81	7.49	7.99	29.4	27.2	2.70	3.35	0	1	
benzylacetone	15	5	10	1.69	0.93	5.93	6.32	30.6	30.3	3.47	4.27	3	8	
biphenyl	24	7	17	2.43	0.55	10.44	10.00	10.2	9.97	2.81	3.35	1		
bromobenzene	21	11	10	3.48	1.09	5.39	7.34	12.6	11.2	4.36	4.63	3	1	
butyric acid ethyl ester	16	5	11	1.65	0.99	3.50	4.64	5.24	7.19	1.99	1.85			
caffeine	25	5	20	1.41	0.78	24.73	18.00	31.5	27.7	5.65	7.08	7	0	
capric acid ethyl ester	16	4	12	3.57	0.56	13.02	13.61	17.0	16.6	0.83	1.65	9	4	
caprylic acid ethyl ester	16	4	12	2.02	0.92	9.87	10.35	16.2	13.7	2.18	1.44	0	1	
chlorobenzene	21	5	16	2.71	0.80	4.57	6.18	8.99	11.0	1.46	4.06		8	
chromium(III) acetylacetonate	104	35	69	3.31	1.41	15.95	15.36	45.6	44.8	7.17	6.69	9	6	
chrysene	4	3	1	4.79	0.87	15.63	17.76	17.6	19.7	n.d.	n.d.	5	3	
citral	15	5	10	2.43	0.67	10.00	7.95	7.57	12.9	5.14	3.99		9	
cobalt(III) acetylacetonate	38	13	25	1.68	0.90	12.06	11.25	47.0	47.3	1.57	2.49	2	2	
copper(II) trifluoroacetylacetonate	12	5	7	6.92	0.60	31.29	41.24	46.0	57.7	10.5	1.90	8	8	3
cycloheptanone	8	3	5	3.91	0.76	25.74	22.98	14.5	21.3	7.41	1.45	6	9	
cyclononanone	8	3	5	3.69	0.82	17.91	17.45	22.7	22.6	2.36	2.78	0	8	

cyclopentanone	8	2	6	4.88	0.24	19.03	20.73	9.43	9.67	1.61	0.81
DHA ethyl ester	65	23	42	0.86	0.41	17.03	17.47	30.9	30.9	1.35	1.44
								8	7		
DHA methyl ester	17	2	15	0.90	0.57	17.80	17.31	25.1	32.5	0.85	0.94
								4	5		
dibenzo-24-crown-8	28	9	19	0.92	0.51	12.30	12.93	52.9	50.6	2.05	1.96
								2	9		
dibenzyl ether	15	5	10	0.88	0.53	3.70	6.13	35.8	39.3	5.35	2.51
								0	4		
diethyl ether	17	2	15	4.76	0.83	12.13	10.33	4.36	8.64	2.39	6.06
diisopropyl ether	15	3	12	7.14	0.58	9.94	6.44	7.66	12.4	4.73	8.47
									7		
diolein	9	5	4	2.55	0.74	23.56	23.85	49.0	48.0	1.74	1.82
								2	3		
Disperse blue 14	47	14	33	3.52	1.29	20.18	20.84	20.1	20.1	3.98	2.17
								8	0		
Disperse orange 11	65	18	47	3.08	1.80	20.01	20.46	14.0	14.9	3.57	3.70
								9	8		
D-limonene	15	4	11	2.38	0.75	10.72	8.80	7.75	7.03	4.80	3.77
docosahexaenoic acid (DHA)	63	22	41	1.18	0.60	9.10	7.42	49.0	47.0	2.12	1.53
								4	6		
eicosapentaenoic acid (EPA)	55	15	40	1.30	0.54	7.23	8.00	41.1	41.8	2.07	1.65
								9	6		
EPA ethyl ester	48	20	28	0.52	0.44	14.74	15.14	30.6	29.4	0.99	1.12
								6	0		
EPA methyl ester	17	6	11	1.40	0.43	17.37	17.36	30.0	30.7	0.51	0.46
								0	4		
ethanol	24	7	17	2.09	0.92	15.75	9.46	8.94	9.98	3.76	2.59
ethyl acetate	16	5	11	3.77	0.47	23.95	14.70	5.22	7.88	6.00	7.99
ethyl benzoate	15	7	8	4.04	0.73	4.43	3.39	21.2	28.0	3.97	1.95
								2	5		

ethylbenzene	15	4	11	0.68	0.69	7.65	7.36	4.89	4.81	2.67	2.31
eugenol	15	3	12	3.05	1.03	18.44	17.00	33.7	41.1	7.21	1.63
								8	7		
ferrocene	107	30	77	1.90	0.79	16.64	17.74	17.5	17.2	6.48	6.88
								6	1		
fluorobenzene	15	3	12	6.71	0.68	11.27	10.98	8.50	11.4	4.38	4.24
									1		
geraniol	4	0	4	n.d.	0.76	n.d.	3.34	n.d.	35.8	n.d.	0.38
									0		
hexachlorobenzene	14	4	10	4.16	0.60	6.75	12.69	14.0	14.2	4.17	4.49
									9	2	
Ibuprofen	99	27	72	1.55	0.73	9.35	10.05	18.9	17.9	4.14	4.13
									1	0	
iodobenzene	20	4	16	2.76	0.61	3.06	10.50	17.3	21.7	3.51	2.72
									1	5	
i-propylbenzene	35	6	30	1.01	0.91	7.19	9.06	7.82	7.49	3.05	2.23
isobutylbenzene	15	8	7	1.99	0.59	3.94	5.63	15.2	19.7	3.20	1.98
									7	1	
L-carvone	27	10	17	1.83	0.88	2.85	4.13	24.2	24.5	2.58	2.69
									5	0	
linalool	15	4	11	2.52	0.77	10.19	6.16	9.72	10.7	3.02	4.42
									7		
linoleic acid	71	27	44	1.09	0.83	9.54	9.68	37.3	38.4	3.54	3.98
									3	1	
linoleic acid methyl ester	20	3	17	2.03	0.77	13.41	15.80	37.3	37.4	1.91	1.53
									0	4	
L-menthone	23	5	18	1.74	1.13	5.67	5.04	21.6	19.2	2.65	2.85
									0	8	
methanol	10	7	3	2.55	0.20	14.11	23.06	18.3	17.8	3.96	0.18
									1	5	
monoolein	11	4	7	0.65	0.69	6.78	9.81	41.6	44.5	1.53	1.26

n-pentylbenzene	31	10	21	3.48	0.51	8.25	8.43	16.4	17.5	4.51	2.58		
												2	3
n-propylbenzene	60	13	47	3.44	0.87	7.41	13.12	9.47	7.64	4.77	4.56		
n-tetradecane	5	1	4	6.13	0.28	31.93	40.49	1.56	14.9	11.5	0.78		
												1	7
n-undecane	5	3	2	6.66	0.50	40.84	41.08	22.1	21.2	2.94	0.00		
												6	9
oleic acid	19	4	15	1.47	0.30	11.96	9.51	40.5	39.8	2.44	2.14		
												7	7
oleic acid ethyl ester	5	2	3	3.16	0.40	9.75	3.04	28.8	29.3	2.74	0.23		
												5	3
oleic acid methyl ester	21	5	16	2.35	0.42	8.94	6.78	29.8	28.2	5.15	3.31		
												2	3
palladium(II) acetylacetonate	125	41	84	1.26	0.62	20.59	22.59	38.4	37.8	5.22	4.62		
												1	8
palmitic acid ethyl ester	17	5	12	1.67	0.53	14.84	15.27	31.3	27.9	0.26	0.75		
												0	0
p-dichlorobenzene	13	4	9	3.73	0.48	10.94	10.46	13.6	18.5	5.69	2.87		
												1	3
phenanthrene	25	7	18	6.74	1.05	14.93	15.69	4.20	5.83	5.04	4.57		
phenol	109	27	82	2.31	0.48	21.40	21.14	9.58	10.6	5.55	4.47		
												2	
phenylacetic acid	16	2	14	1.74	0.62	5.82	4.02	17.6	15.5	1.83	1.86		
												9	8
phenylacetylene	15	5	10	1.49	0.72	6.18	8.61	16.4	15.8	1.69	1.58		
												5	4
phenylbutazone	78	20	58	1.78	0.69	9.65	7.67	34.1	32.1	5.29	5.98		
												3	8
phenylmethanol	15	4	11	1.51	0.53	14.14	13.83	22.2	22.7	3.28	2.26		
												3	8
p-xylene	7	1	6	0.91	2.19	7.05	6.46	1.95	4.95	1.14	4.02		

pyrene	21	9	12	2.99	1.08	8.64	10.58	24.5	16.7	3.74	3.32	7	2
sec-butylbenzene	15	4	11	1.63	0.65	3.86	4.03	22.7	18.9	4.68	2.83	3	5
squalene	5	1	4	2.35	1.30	33.70	29.88	18.6	26.7	3.93	1.15	7	5
stearic acid	4	2	2	2.09	1.00	50.13	50.51	20.7	21.4	0.77	0.00	1	3
stearic acid ethyl ester	17	5	12	0.62	0.59	24.10	24.53	29.6	24.0	1.27	0.99	4	6
styrene	15	6	9	4.28	0.68	6.82	4.43	12.8	15.9	5.13	3.94	9	4
tert-butylbenzene	15	5	10	1.53	0.75	8.26	7.74	15.6	13.2	3.33	3.95	6	1
tetrahydrofuran	15	2	13	4.36	0.30	8.08	17.18	2.38	12.2	12.2	4.50	1	9
thenoyltrifluoroacetone	15	2	13	1.88	1.13	35.57	29.37	44.6	48.0	2.75	3.43	6	9
toluene	41	10	31	3.26	0.68	10.33	8.44	4.57	4.65	9.01	8.26		
triarachidonin	27	8	19	1.56	0.82	18.29	17.15	70.0	70.3	0.67	0.88	2	6
trierucin	101	30	71	2.22	0.67	15.40	12.80	80.4	83.1	2.68	3.06	5	0
trifluoroacetylacetone	15	6	9	3.62	0.36	2.88	4.30	12.2	10.4	3.14	1.46	8	3
trinervonin	38	13	25	1.54	0.60	16.15	16.86	84.2	82.6	2.84	2.97	5	6
triolein	14	6	8	2.11	0.44	30.59	28.66	54.8	53.4	3.95	2.57	0	0
ubiquinone CoQ10	80	27	53	2.22	0.84	12.49	14.61	71.6	71.4	4.00	4.21	8	8

vanillin	15	5	10	2.36	1.40	11.45	13.23	23.7	25.0	2.18	2.08
								9	5		
vitamin K1	17	8	9	2.37	1.93	25.75	28.57	31.3	31.1	2.45	2.56
								0	2		
vitamin K3	22	4	18	3.05	1.27	8.83	9.75	12.6	11.8	4.75	2.56
								4	9		
water	24	3	21	5.73	0.25	54.04	56.51	21.5	12.0	3.21	4.32
								7	7		

n.d.: not determined.

Figure 4: Experimental and calculated diffusivities (ML-GB-SC-CO₂ model) in terms of free-volume theory coordinates for (a) acetone and (b) 1-methylnaphthalene.

Figure 5: Experimental and calculated diffusivities (ML-GB-SC-CO₂ model) in terms of Stokes-Einstein coordinates for (a) acetone and (b) 1-methylnaphthalene.

4 Conclusions

In this work a machine learning model for the prediction of binary diffusivities in SC-CO₂ was developed. This model was trained and validated by splitting a database containing 13 properties of 174 systems and 4917 points into training and test sets. Several learning algorithms were tested (Multilinear Regression, k -Nearest Neighbors, Decision Tree, Random Forest and Gradient Boosted). The best results were found using the Gradient Boosted algorithm, which presented an average deviation of 2.58 % for the test set. This model takes five input properties/variables which are readily available for multiple solutes: temperature, solvent density, solute molar mass, solute critical pressure and solute acentric factor. Results were compared with the classic diffusivity equations of Wilke-Chang, Lai-Tan, and Dymond-Hildebrand-Batschinski, which demonstrated worse performance for the same data with deviations of 12.41 %, 26.01 % and 4.27 %, respectively. Although the Dymond-Hildebrand-Batschinski model shows similar performance, it requires *a priori* experimental data to fit the system parameters, which is not always possible.

Acknowledgements

This work was developed within the scope of the project CICECO-Aveiro Institute of Materials, UIDB/50011/2020 & UIDP/50011/2020, financed by national funds through the Foundation for Science and Technology/MCTES, as well as the Multibiorefinery project (POCI-01-0145-FEDER-016403). Bruno Zêzere thanks FCT for PhD grant SFRH/BD/137751/2018.

References

- [1] P C Wankat. *Rate-controlled separations*. Blackie Academic & Professional, Great Yarmouth, 1994.
- [2] Eduardo L G Oliveira, Armando J D Silvestre, and Carlos M. Silva. Review of kinetic models for supercritical fluid extraction. *Chemical Engineering Research & Design*, 89(7A):1104–1117, 2011.
- [3] James J. Carberry. *Chemical and catalytic reaction engineering*. McGraw-Hill, 1971.
- [4] Denis Prat, John Hayler, and Andy Wells. A survey of solvent selection guides. *Green Chemistry*, 16(10):4546–4551, oct 2014.
- [5] M M R de Melo, A J D Silvestre, and C M Silva. Supercritical fluid extraction of vegetable matrices: applications, trends and future perspectives of a convincing green technology. *The Journal of Supercritical Fluids*, 2014.
- [6] Bruno Zêzere, Ana L. Magalhães, Inês Portugal, and Carlos Manuel Silva. Diffusion coefficients of eucalyptol at infinite dilution in compressed liquid ethanol and in supercritical CO₂ ethanol mixtures. *Journal of Supercritical Fluids*, 133:297–308, mar 2018.
- [7] C. R. Wilke and Pin Chang. Correlation of diffusion coefficients in dilute solutions. *AIChE Journal*, 1(2):264–270, jun 1955.
- [8] Ching Chih Lai and Chung Sung Tan. Measurement of molecular diffusion coefficients in supercritical carbon dioxide using a coated capillary column. *Industrial and Engineering Chemistry Research*, 1995.
- [9] Raquel V. Vaz, Ana L. Magalhães, and Carlos M. Silva. Improved hydrodynamic equations for the accurate prediction of diffusivities in supercritical carbon dioxide. *Fluid Phase Equilibria*, 2013.
- [10] Raquel V. Vaz, Ana L. Magalhães, and Carlos M. Silva. Improved Stokes-Einstein based

- models for diffusivities in supercritical CO₂. *Journal of the Taiwan Institute of Chemical Engineers*, 45(4):1280–1284, 2014.
- [11] J. H. Dymond. Corrected Enskog theory and the transport coefficients of liquids. *The Journal of Chemical Physics*, 1974.
- [12] J. H. Dymond, E. Bich, E. Vogel, W. A. Wakeham, V. Vesovic, and M. J. Assael. Dense Fluids. In Jürgen Millat, J. H. Dymond, and C. A. Nieto de Castro, editors, *Transport Properties of Fluids*, pages 66–112. Cambridge University Press, 1996.
- [13] C.M. Silva and H. Liu. Modelling of Transport Properties of Hard Sphere Fluids and Related Systems, and its Applications. In *Theory and Simulation of Hard-Sphere Fluids and Related Systems*, pages 383–492. Springer, Berlin, 2008.
- [14] F. Gharagheizi and M. Sattari. Estimation of molecular diffusivity of pure chemicals in water: A quantitative structure-property relationship study. *SAR and QSAR in Environmental Research*, 2009.
- [15] Aboozar Khajeh and Mohammad Reza Raeesi. Diffusion coefficient prediction of acids in water at infinite dilution by QSPR method. *Structural Chemistry*, 2012.
- [16] Raquel V. Vaz, Ana L. Magalhães, Daniel L.A. Fernandes, and Carlos M. Silva. Universal correlation of self-diffusion coefficients of model and real fluids based on residual entropy scaling law. *Chemical Engineering Science*, 2012.
- [17] Reza Beigzadeh, Masoud Mahini, and Seyed Reza Shabani. Developing a feed forward neural network multilayer model for prediction of binary diffusion coefficient in liquids. *Fluid Phase Equilibria*, 2012.
- [18] R. Eslamloueyan and M. H. Khademi. A neural network-based method for estimation of binary gas diffusivity. *Chemometrics and Intelligent Laboratory Systems*, 2010.
- [19] Seyyed Alireza Mirkhani, Farhad Gharagheizi, and Mehdi Sattari. A QSPR model for prediction of diffusion coefficient of non-electrolyte organic compounds in air at ambient condition. *Chemosphere*, 2012.
- [20] Alireza Abbasi and Reza Eslamloueyan. Determination of binary diffusion coefficients of hydrocarbon mixtures using MLP and ANFIS networks based on QSPR method. *Chemometrics and Intelligent Laboratory Systems*, 2014.
- [21] Kenneth S. Pitzer and Donald R. Schreiber. Improving equation-of-state accuracy in the critical region; equations for carbon dioxide and neopentane as examples. *Fluid Phase*

- Equilibria*, 1988.
- [22] V. V. Altunin and M. A. Sakhabetdinov. Viscosity of liquid and gaseous carbon dioxide at temperatures of 220 - 1300 K and pressures up to 1200 bar. *Thermal Engineering (English translation of Teploenergetika)*, 1972.
- [23] Robert C. Reid, John M. Prausnitz, and Bruce E. Poling. *The Properties of Gases & Liquids*. McGraw-Hill, 4th edition, 1987.
- [24] K.G. Joback. *A Unified Approach to Physical Property Estimation Using Multivariate Statistical Techniques*. Master thesis, Massachusetts Institute of Technology, Cambridge, MA, 1984.
- [25] K. G. Joback and R. C. Reid. Estimation of pure-component properties from group-contributions. *Chemical Engineering Communications*, 57(1-6):233–243, jul 1987.
- [26] G. Raam Somayajulu. Estimation Procedures for Critical Constants. *Journal of Chemical and Engineering Data*, 1989.
- [27] Ana L. Magalhães, Francisco A. Da Silva, and Carlos M. Silva. Free-volume model for the diffusion coefficients of solutes at infinite dilution in supercritical CO₂ and liquid H₂O. *Journal of Supercritical Fluids*, 74:92–104, 2013.
- [28] K. M. Klinecicz and R. C. Reid. Estimation of critical properties with group contribution methods. *AIChE Journal*, 1984.
- [29] Robert H. Perry and Don W. Green. *Perry's Chemical Engineers' Handbook*, volume Sixth edit. McGraw-Hill, 6th edition, 1997.
- [30] Leonidas Constantinou and Rafiqul Gani. New group contribution method for estimating properties of pure compounds. *AIChE Journal*, 40(10):1697–1710, oct 1994.
- [31] Byung Ik Lee and Michael G. Kesler. A generalized thermodynamic correlation based on three-parameter corresponding states. *AIChE Journal*, 1975.
- [32] Kenneth S. Pitzer, David Z. Lippmann, R. F. Curl, Charles M. Huggins, and Donald E. Petersen. The Volumetric and Thermodynamic Properties of Fluids. II. Compressibility Factor, Vapor Pressure and Entropy of Vaporization. *Journal of the American Chemical Society*, 1955.
- [33] Hongqin Liu, Carlos M. Silva, and Eugénia A. Macedo. New Equations for Tracer Diffusion Coefficients of Solutes in Supercritical and Liquid Solvents Based on the Lennard-Jones Fluid Model. *Industrial & Engineering Chemistry Research*,

- 36(1):246–252, 1997.
- [34] Andriy Burkov. *The Hundred-Page Machine Learning Book*. Andriy Burkov, 2019.
- [35] Fabian Pedregosa, Gael Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, oct 2011.
- [36] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2nd edition, 2009.
- [37] Andreas C. Müller and Sarah Guido. *Introduction to Machine Learning with Python: a guide for data scientists*. O’Reilly Media, Inc., 2016.
- [38] B E Poling, J M Prausnitz, and J P O’Connell. *The Properties of Gases and Liquids*. McGraw-Hill, Singapore, 5th edition, 2000.

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:

Journal Pre-proof

University of Aveiro, November 30th 2020

Dear Professor Stratos Pistikopoulos,
Editor-in-Chief of the Journal of Molecular Liquids

On behalf of all co-authors I hereby confirm that we have no conflict of interests with any researcher or entity, and that our work is being uniquely submitted to the *Journal of Molecular Liquids*.

Carlos Manuel Silva

Chem. Eng. Group | CICECO - Aveiro Institute of Materials
Associate Professor | Department of Chemistry | University of Aveiro
Campus Universitário de Santiago
3810-193 Aveiro | Portugal

E-mail: carlos.manuel@ua.pt | Web: www.eqjchem.com
Phone: +351 234 401 549 | Ext: 24928

Highlights

- New predictive model to estimate diffusivities in supercritical carbon dioxide.
- The new machine learning model was trained with a database of 174 binary systems.
- It was compared with several classical models, such as the Wilke-Chang equation.
- The machine learning model provided the best performance with errors of 2.58 %.

Journal Pre-proof