1	A Three-Pronged Computational Approach for Evaluating Density Based
2	Semi Empirical Equations of Supercritical Extraction Process and Data
3	
4	Srinidhi <sup>1, *</sup>
5	
6	<sup>1</sup> Department of Chemical and Biological Engineering, School of Engineering and Applied
7	Sciences, The State University of New York at Buffalo, Buffalo, NY 14260.
8	
9	
10	*Corresponding author, e-mail: <u>srinidh2@buffalo.edu</u> ; <u>msrinidhi2@gmail.com</u>
11	ORCID: <u>0000-0002-5318-8639</u> ;
12	
13	

#### Abstract

Software programs for parameter estimation, phase visualization and predictive 16 modeling of supercritical extraction process and data using algorithms is presented in this work. 17 A contextually appropriate, iterative, ordinary least squares estimation and selection method is 18 developed for estimating model coefficients of density based semi empirical model equations 19 associated with this process and data. Visualization of the phase behaviors projected by the 20 specific density based semiempirical model equation(s) is also performed iteratively by plotting 21 three-dimensional surfaces involving the state variables and solute solubility mole fraction. 22 Predictive modeling of input empirical data has been implemented using three supervised 23 machine learning algorithms (Multilayer perceptron, K-nearest neighbors and Support vector 24 regression). Hyperparameter optimization of the machine learning algorithms is performed 25 prior to prediction. Detailed analysis of the prediction is conducted by using standard scoring 26 metrics and descriptive charts. Theoretical inference and discrepancies regarding the predicted 27 window of maximum/optimal solubility, modeling efficiency, vapor liquid equilibrium and 28 phase behaviors projected by the model equations have been elucidated from the program 29 outputs. In summary, these programs are unique, accurate, reliable and simple computational 30 tools for evaluating/designing density based semiempirical equation(s) of supercritical 31 extraction process and associated data. 32

33

14

15

Keywords: Parameter Estimation, Phase Visualization, Predictive modeling, Ordinary Least
Squares, Machine Learning.

- 36
- 37

# Introduction

38

Theoretical, Empirical and Semi empirical Models are being developed and studied for modeling and understanding Super/subcritical fluid extraction processes (Huang et al. 2012; Rai et al. 2014). In particular, Density based Semiempirical model equations (DBSE Model Equations) are very popular and are being designed for modeling this process and therefore is part of a growing body of research (Hawthorne 1990; Herrero et al. 2010; Knez et al. 2013; Alwi and Garlapati 2021a).

<sup>45</sup> Novel DBSE models are developed with the aim to capture (approximate) and reproduce data
<sup>46</sup> specific non linearity and complexity (dynamic and non-dynamic behavior) in the process.
<sup>47</sup> Modeling in this scenario is primarily focused on the operating range of the process parameters

observed during desired output/yield levels (Tabernero et al. 2010). Unfortunately, in most 48 cases, this window (presumably rich in information) is narrow and is solute/process specific. 49 Almost every study describing a novel DBSE model have proceeded by distilling facts about 50 the variation in solvating power observed in the process and drawing fundamental relations 51 (from similar studies) between the operating process parameters and the dependent variable 52 [ln(y), T, P, D]. A good and elegant example for this is the study and model presented by 53 (Asgarpour Khansary et al. 2015). Least squares modeling is a subclass of Black box modeling 54 and has been extensively employed for estimating model parameters, their confidence regions 55 (Bounds/Intervals) and importantly for identifying causation of variance in linear models. 56 Herein, Ordinary least squares estimation method is used for estimating parameter coefficients 57 (and their confidence regions) present in DBSE Model equations (Lakshmi et al. 2021). 58

Further, A necessary requirement for the design of DBSE models is the qualitative and 59 quantitative knowledge of phase behavior of components in the reaction mix during the process. 60 Phase diagrams illustrate important differentials in vapor pressure curves of pure CO2 and other 61 reaction components in the presence of solutes. This information is crucial for accurately 62 identifying operating conditions wherein melting of the reaction mix leading to a desirable 63 solute rich liquid phase occurs. In essence, phase diagrams are central to the process of finding 64 regions (boundaries) of importance in the P-T-D-x (Pressure-Temperature-Density-solute 65 solubility mole fraction) projections, wherein separations and extraction is actually possible 66 (and feasible) and occurs in reality (Bartle et al. 1991). These regions (phase/parameter 67 boundaries) depict equilibrium planes and latency of reaction mix that aid in process design and 68 this is considered as a multifaceted and multi-attribute dependent endeavour. These attributes 69 can be and are not limited to, 70

Regions where solvent compression occurs leading to repulsive solute-solvent
 interactions causing undesired immiscibility.

Regions where two-phase retrograde condensation/crystallization occurs near the lower
 and upper crossover regions/planes/edges.

Regions (edges/paths/points/trajectories) depicting the component(s) latency (phase
 change), chemical potential thermal stability of the solute leading to variations in
 solvating power/effect. Physical properties of solutes vary widely and significantly
 amount to differences during solute solubility prediction.

79 Machine learning algorithms, in recent years, are gaining importance and are being developed 80 for predictive modeling for engineering applications. ML algorithms can accommodate

(consider) 'n' number of parameters, and therefore can predictively model processes with 81 desired tolerance, precision and accuracy. Invaluable for accountability and research 82 applications, hyperparameters associated with ML algorithms offers the choice of model 83 optimization and validation. Standardized ML algorithms are applied to model a multitude of 84 phenomena/processes in Engineering (Selvaratnam and Koodali 2021). Therefore, with the fast 85 parametrization and modeling of analytical and industrial processes, supervised learning 86 models like, Regression, Multilayer Perceptron, Support Vector Machine and K-nearest 87 neighbours are (can also be) specially applied to these processes. For Chemistry and Chemical 88 Engineering applications, A number of Software program packages based on supervised 89 learning are already available and are always under continuous development (Khatib and de 90 Jong 2020). In recent years, estimating/predicting solute solubility during the supercritical fluid 91 extraction is gaining importance and necessitates predictive modeling of this process (Butler et 92 al. 2018; Schweidtmann et al. 2021; Roach et al. 2023). The reliable and utilitarian software 93 program can be used to accurately describe extant pattern and behavior in the measured data 94 associated with this process and possibly beyond the regions and scope of this measured 95 empirical data for reaching higher levels of process interpretation and accurate predictive 96 capabilities. With this as the goal, the predictive modeling program described here has been 97 written and focused to meet this expectation(s). Further, the complete work (workflow) 98 presented here, is also designed for visualization and for explicating the phase behavior of 99 existing (and newer) model equations and for evaluating the boundedness of the estimated 100parameter space. This workflow is holistic and is particularly useful for designing newer, 101 efficient (accurate and precise) equations heuristically. Conveniently, As previously mentioned, 102 a one-time-run-all code has been provided for implementing state-of-the-art machine learning 103 algorithms for predictive modeling of DBSE model equation associated data. When correctly 104 deployed, this work could potentially reach the helm of this growing body of research from this 105 three-pronged computational modeling approach. To summarize, the programs are stand alone, 106 simple, unique, computationally economic and are also easy to implement. The objectives and 107 Software being postured here in this article are listed below, 108

A MATLAB program for estimating and comparatively analysing, parameters of
 extant/newly developed density based semi empirical model equations of supercritical
 fluid extraction process comprising of variables [ln(y), T, P, D] using ordinary least
 squares parameter estimation method.

A MATLAB program for visualizing parameter profiles and Phase behaviors of DBSE
 model equations using 3D surface plots.

A Python based Jupyter Notebook for implementing supervised machine learning algorithms (Multilayer Perceptron, K nearest neighbours and Support vector machines) based on experimental data involving the variables (Temperature (T), Pressure (P), Density (D) and Solute solubility Mole fraction (y)).
 Provide concluding remarks about the program scripts, its usage and availability.
 Experimental

## Description of Data: Input Matrices and Parameter Description

122 123

121

The MATLAB (Matlab 1984) and Python program scripts presented in this work requires two input matrices. First, Consider, the Input data as a matrix where in,  $Data \in \mathbb{R}^n$  and  $n \in \mathbb{Z}$ , then,

127 
$$Data_{i,4} = \begin{bmatrix} T_{1,1} & P_{1,2} & D_{1,3} & y_{1,4} \\ \vdots & \vdots & \vdots & \vdots \\ T_{i,1} & P_{i,2} & D_{i,3} & y_{i,4} \end{bmatrix}$$
 (1)

Where T is temperature in Kelvin, P is pressure in Mpa, D is density in Kg/m3 and y is solubility mole fraction of the solute in the reaction mix. The index 'i', runs over the entire column of a single feature. This is the first input data matrix required and is parsed by the scripts via the Input\_Data.xlsx file.



Fig. 1 Flow chart illustrating a single iteration by the parameter estimation, 3D visualizationand Predictive modeling program scripts



135

Fig. 2 (a) Heat Map plot of correlation values of input parameters (Temperature, Pressure,
Density and Solute Mole fraction). (b) Parameter pair plot of data points including all
combinations of input parameters for illustrating patterns present among variable pairs

139

The second (required) matrix is comprised of the terms of the input density based semi
empirical equations. For illustration, consider a simple four parameter (however, users can input
any number of terms) linear model equation and its basic generalization,

143  $\ln(y) = A + B[T] + C[P] + D[\rho] \equiv Y = p_1[Term1] + p_2[Term2] + p_3[Term3] + p_4[Term4]$  (2) 144 Where, A, B, C, D corresponds to p1, p2, p3, p4 and are the parameter coefficients or estimands 145 of the DBSE model equation given above (however, users can input any number of terms). Let

these parameter coefficients be grouped into vector 'P'. Let the terms of the model [term1, 146 term2, term3, term4] be grouped into a vector named as 'Terms'. For the estimation of the 147 model coefficients in [P] and for obtaining parameter estimates  $\hat{p}$ , the terms of the sampled 148 DBSE model equations are input into respective cells of rows particular to each model equation 149 in a separate file (Models\_Equations.xlsx). These are the two input matrices required by the 150 MATLAB based parameter estimation script and the visualization script. A toy input data 151 sample containing 1000 experiments along with a sample of ten randomly selected, 152 semiempirical equations have been used for producing the output present in this article. The 153 possible modification path traversed by the data in a single iteration is illustrated (Fig. 1). Also, 154 the Input data is initially analyzed using basic statistical metrics in the Jupyter Notebook and 155 the outputs (Correlation heat map and parameter pair plot) are depicted (Fig. 2 a, b). Refer to 156 the user guide (given in the repository) for information on using these program scripts for 157 custom data and model equations (existing/newly proposed). The user guide also provides 158 information regarding the preselection of the base model along with the descriptions of the 159 randomly sampled model equations present in the unmodified file (Models\_Equations.xlsx). 160

- 161
- 162

## Parameter Estimation: Ordinary Least Squares Method

163

Estimation of parameter coefficients represented in the vector P is performed using the method of Ordinary Least Squares Parameter Estimation (Dismuke and C R Lindrooth 2006) in the MATLAB program script (DBSE\_OLS\_Estimation.m). A concise development of the implemented algorithm is presented. Consider a representation of a DBSE model equation in the form of the classical linear regression model,

169 
$$Y_{i,1} = [Terms]_{i,k}[P] + \varepsilon_{i,1}$$
 (3)

170 Let the assumptions, about the error in the models be, errors are additive, uncorrelated, has zero171 mean and has constant variance.

173 
$$E(\varepsilon\varepsilon^{1}) = \sigma^{2}I_{i}$$
(4)

174 Where  $\varepsilon$  is the residual vector and  $\sigma^2$  is the variance of the residual. Further, let the data 175 substituted, matrix of terms 'Terms' be represented for brevity as X and let Y be the vector of 176 natural logarithm of solubility mole fraction values. Then the ordinary least squares estimator 177  $\hat{p}$  is given by,

178 
$$\hat{p} = [X^T X]^{-1} X^T Y$$
 (5)

179 The vector of residuals  $\varepsilon$  is given by,

$$180 \ \varepsilon = Y - X\hat{p} \tag{6}$$

The confidence intervals (bounds) of the estimates are computed at 95% confidence level. Further, model selection is iteratively performed using an F-Statistic score (Belitser et al. 2011) for each model equation relative to a preselected base model (This is input in the first row of the Models\_Equations.xlsx file). Let the residual sum of squares for the DBSE model of a particular iteration and the same for base model be,

$$R_{ols}^{model} = \varepsilon_{ols}^T \varepsilon R_{ols}^{base} = \varepsilon_{base}^T \varepsilon$$
(7)

187 Then the equation for an F-score metric-based model selection is,

$$\frac{\binom{R_{ols}^{base} - R_{ols}^{model}}{\binom{n_{p,0} - n_{p,base}}{(n_{p,0} - n_{p,base})}} > F_{(n_{p,0} - n_{p,base}),(n - n_{p,0})}^{0.05}$$
(8)

Where n<sub>p,0</sub> is the number of parameters in the current iteration and n is the number of data points 189 (experiments) in the parsed input data and  $n_{p,base}$  is the number of parameters in the base model. 190 In the data driven paradigm where modeling is focused on fitting a specific sample of empirical 191 data, this automated selection procedure is beneficial for decimating lower quality equations 192 and for identifying the most contextually appropriate one(s). Further, error metrics namely, 193 mean squared error (MSE), Root Mean Squared Error (RMSE), Mean Absolute error (MAE) 194 and Percentage Absolute Average Relative Deviation (% AARD) were computed between 195 experimental and predicted solubility using the expressions, 196

197 Mean Squared Error = 
$$\frac{1}{n} \sum_{i=1}^{n} \left( \ln \left( y \right)_{i}^{pred} - \ln \left( y \right)_{i}^{exp} \right)^{2}$$
(9)

198 Root Mean Squared Error = 
$$\sqrt{\frac{1}{n}\sum_{i=1}^{n} \left(\ln\left(y\right)_{i}^{pred} - \ln\left(y\right)_{i}^{exp}\right)^{2}}$$
 (10)

199 Mean Absolute Error 
$$= \frac{1}{n} \sum_{i=1}^{n} \left| \left( \ln(y)_i^{pred} - \ln(y)_i^{exp} \right) \right|$$
 (11)

200 %AARD = 
$$\frac{100}{n} \sum_{i=1}^{n} \frac{\left|\ln(y)_{i}^{pred} - \ln(y)_{i}^{exp}\right|}{\ln(y)_{i}^{exp}}$$
 (12)

Error metrics have been computed using natural logarithm of solubility mole fraction values
for predictions after parameter estimation and actual solubility mole fraction values have been
used for predictions from predictive modeling.

204

## 205

#### Visualization of Phase Behaviour Projected by DBSE model Equations:

206

<sup>207</sup> Visualization of Phase behavior using three dimensional surfaces of the input DBSE model <sup>208</sup> equation is also implemented using MATLAB. The MATLAB script (DBSE\_3D\_Viewer.m), 209 requires, model equations and empirical data (Input\_Data.xlsx and Models\_Equations.xlsx)
210 along with the estimates (Parameter\_Predictions\_Results.xlsx) and iteratively plots three
211 dimensional surfaces of the model equations using finitely spaced grid points of the parameters
212 present in the particular DBSE model equation in the iteration.

Three surfaces are plotted by this script namely, Pressure-Temperature-Solute mole 213 fraction, Density-Pressure-Solute mole fraction and, Density-Temperature-Solute mole 214 fraction. Standard, inbuilt commands from MATLAB are used for plotting the surfaces for all 215 of the input DBSE model equations. The output images are also in the standard interactive 216 MATLAB plot window which allows for altering values of axes to obtain surfaces (Rovenski 217 2010). Notedly, empirical data is used by this script only for finalizing extreme values of the 218 grid points used for plotting these surfaces. Therefore, the surfaces plotted by this script 219 illustrate phase behavior and vapor liquid equilibrium data projected by the specific DBSE 220 model equation and these surfaces are not influenced by the pattern prevalent in the input 221 empirical data. Finally, this script exports all three surfaces plotted for a DBSE model equation 222 as subplots in a single image (.jpg) format. 223

- 224
- 225

#### Prediction of Solute Solubility: Machine Learning Algorithms

226

Three Supervised Machine learning algorithms have been implemented using the Python 227 module, Sklearn (Pedregosa et al. 2011) in a single Jupyter notebook 228 (DBSE\_Predictive\_Modeling.ipynb) (Menke 2020). This Notebook, using input empirical data, 229 in a single run, implements the Multilayer perceptron, K-nearest Neighbours regression and 230 Support Vector regression algorithms before performing detailed and comparative analysis on 231 the predictions and results. Standardized metrics are used for performing validation and analysis 232 of results. Numpy (Oliphant 2006), Openpyxl, Pandas (W McKinney 2011), Matplotlib (Hunter 233 2007) are among the python packages used for implementing these algorithms. This script 234 requires empirical data (experiments in rows complete with Pressure, Temperature, Density and 235 the resultant, solute mole fraction) characteristic to density based semi empirical model 236 equations. Also, the input parameter space is not exhaustive and can incorporate additional 237 parameters based on preference. Descriptions of the implemented algorithms and their tuneable 238 hyperparameters are provided in the subsequent paragraphs. 239

- 240
- 241

#### Multilayer Perceptron Regression [MLP]

243

Multilayer Perceptron [MLP] is a fully connected class of feed forward artificial neural 244 networks classified as a supervised machine learning algorithm. This framework consists of 245 updatable, weight assigned nodes called neurons that are sorted into three types of fully 246 connected layers namely, input layer, hidden layer(s) and an output layer. During the training 247 of a single instance (experiment), parameter (feature) information is fed into the input layer 248 which is then transmitted to the next hidden layer(s) where activation function(s) modify this 249 information for final modification in the output layer. The output layer, using an activation 250 function, modifies the received information and provides data output. This output is the 251 prediction value of the algorithm. Information modification during training (learning) results in 252 the updation of the initialized weights (associated with neurons and connections) from the 253 previous learning iteration (Murtagh 1991). In MLP, for obtaining accurate and precise output 254 (solute solubility mole fraction), hyperparameter search space for size of hidden layer, neurons, 255 activation functions, learning rate, data split ratio, solver, alpha value etc can be easily 256 optimized in the notebook based on preference and data. Theoretical explanation and 257 development of the MLP algorithm can be obtained in literature elsewhere (Schilling et al. 258 2015). The results and analysis from this program code are finally saved (Ml\_Results.xlsx). 259

- 260
- 261

## K-Nearest Neighbours Regression [KNN]

262

K- Nearest Neighbours algorithm is a non-parametric, supervised machine learning algorithm. 263 For regression problems, the algorithm learns to predict the target class value based on the k 264 closest training examples (instances or experiments) in the input data. The model during 265 learning (training), performs search in the data pattern space for the closest number of training 266 instances. The results from this search which are the closest 'k' number of training instances 267 (neighbours), are averaged to obtain the prediction value (solute solubility mole fraction) during 268 testing (Kramer 2013). The adjustable/tuneable hyperparameters for this algorithm is the 'k' 269 value (sampling metric) and the distance (closeness) measurement metric (Cunningham and 270 Delany 2022). Here, Euclidean distances are calculated to measure closeness for the 271 preassigned k value which is used to obtain a detailed, comparative, analysis of the prediction 272 which also is saved (Ml\_Results.xlsx). 273

#### Support Vector Regression [SVR]

The support vector regression algorithm is a class of support vector machine algorithm and is 277 also a supervised machine learning algorithm. In fewer sentences, support vector regression 278 algorithm, using a kernel function, tries to map the input parameter variable data to a feature 279 space (usually of higher dimension) and with the aim of minimizing prediction error, tries to 280 find a hyperplane in this feature (parameter) space that maximizes the distance margin between 281 this plane and the closest data points. Theoretical development of the SVR technique and the 282 mechanism behind its prediction capabilities can be obtained in detail here (Smola and 283 Schölkopf 2004). The tuneable hyperparameters here are the kernel function, gamma value and 284 the test-train data split ratio. Scaling of the parameter data has not been implemented for SVR 285 as the pattern present in the parameter space is highly relevant for accurate prediction 286 (Tsirikoglou et al. 2017). The jupyter notebook, after implementing support vector regression, 287 separately provides results which also is saved (Ml\_Results.xlsx). 288

289





#### **Results and discussion**



1200

1200

1211

00





Fig. 3 Standard output (enlarged) for the model equation(s) being iterated from the MATLAB based parameter estimation script. (a) Plot of Experimental (black) v/s Predicted (red) values of the natural logarithm values of solute solubility molefraction. (b) Plot containing normality plots and residual plots for base model equation of choice and the model equation being iterated.
(c) Bar plots pertaining to error metrics for all input equations.

As previously derived, A customized Ordinary least squares estimation method has been 298 implemented to obtain parameter estimates of model equation constants along with confidence 299 intervals in a 'one model equation at a time' iterative rule fashion. This ensures that the 300 parameter (model coefficients) estimates are from a standardized and popularly used method 301 used on all model equations in the batch sample (input using an .xlsx file). Confidence intervals 302 (upper and lower bounds) are estimated for each estimate at 95 percent confidence. 303 Conveniently, the results are saved and exported to retrievable file formats. The pictorial output 304 from this script is shown in (Fig. 3 a - c). Natural logarithm values of solute solubility mole 305 fractions are plotted against number of experiments for both empirical data and predictions 306 made using the estimates (model constants) and state variables (Pressure, temperature and 307 Density) associated with the model equations. Normality plots and residuals of the base model 308 and the model equation (being iteratively estimated) are also charted for ascertaining the nature 309 of the data. The normality and residual plots are shown (Fig. 3b). Normality plots reaffirm the 310 considered assumptions about the residuals while estimating parameter coefficients (Model 311 constants). This step makes sure the estimates are contingent with the assumptions made 312 regarding the data and by extension, also the residuals. In the Fig. 3 b above, the data appear to 313 lie on the line of reference demonstrating the degree of normality present in the sample data. 314 Unfortunately, the large amount of data (from the toy data sample) in the shown residuals plot 315 indicate a pattern and masks the randomly distributed points in the region of interest. This 316 region of interest corresponds to the operating conditions where solute solubility is supposedly 317 maximum/optimum (window of maximum solubility). However, this also will change when 318 different empirical data is used. Scores computed from F Distribution, provide clear, statistical 319 comparison between the model equation being iteratively estimated and the base model 320 equation of choice (Input in the first row in the Models\_Equations.xlsx file). Additionally, 321 excellent inference can be made based on published literature regarding the estimates and 322 selection output produced by this program (Garlapati and Madras 2010; Reddy and Madras 323 2011; Bian et al. 2016; Alwi and Garlapati 2021b). The pictorial illustration indicates the 324 plotting constraints (maximum number of subplots in the image output) associated with the 325 presented code and it is encouraged to consider this factor while sampling model equations. 326 Plotting natural logarithm values of the predicted data against actual solute solubility mole 327 fraction values of the predicted data (from model equations), provides clear distinction and 328 higher resolution of model fit and deviation from empirical data. Errors and residuals are also 329 calculated using natural logarithm values for this important reason. In reality, based on the toy 330

331 sample empirical data, the error metrics and residuals appear to be significantly (desirably) low

when actual solubility values are used as opposed to their natural logarithm values. Mean 332 squared error (MSE), Root Mean Squared Error (RMSE), Mean Absolute error (MAE) and 333 Percentage Absolute Average Relative Deviation (% AARD) values are computed using Eq. 334 (9)-(12), plotted and presented in the form of bar graphs in a single image format (Fig. 3 c). 335 Errors of all model equations appear to only slightly differ indicating superior quality of the 336 sampled toy data. However, as previously mentioned, this too will differ for other empirical 337 data. Due to constraints for assessing and visualizing higher numbers of equations, sampling 338 (ten to fifteen equations) and selection of model equations (for achieving column rank) must be 339 of provided higher quality. However, the code for batch estimation 340 (DBSE\_OLS\_Estimation\_Batch.m), has full capability to estimate as many as a hundred DBSE 341 model equations in a single implementation. 342

In summary, this program script provides parameter estimates of model(s) coefficients along with their confidence regions (intervals). Further, the model selection and identification routine is also favorable for comparative assessment and selection of the best performing model equation all of which are then exported to popular file formats.

## Visualization of Phase Behavior projected by DBSE Model Equations:



**Fig.** 4 Three dimensional surfaces of ln(y)-P-T, ln(y)-D-P, ln(y)-D-T. (a) This plot is the only standard output produced by the MATLAB based visualization script. (b) Two dimensional, color coded contour plot of P-T, P-D, D-T obtained from the same MATLAB interactive plot window. The projections for these plots are visible on the respective 3D surface (a). (c) Two dimensional, color coded contour plot of ln(y)-T, ln(y)-D, ln(y)-D obtained from the MATLAB interactive plot window.

356 The three-dimensional surfaces of the P-T-D state variables and the natural logarithm values of solute solubility mole fraction obtained from this script for visualization is illustrated in Fig. 4 357 a - c. The interactive nature of the MATLAB surface plot window and the ease with which axes 358 values of the plotted surface can be altered makes the obtained pictorial output invaluable for 359 evaluating the phase equilibria characteristic to the respective DBSE model equation. Fig. 4 a 360 shows a grab of the three surfaces [P-T-ln(y), P-D-ln(y), T-D-ln(y)] arranged as subplots from 361 a single interactive (image) window output. As previously mentioned, Grabs of two-362 dimensional plots (Fig. 4 b - c) can be obtained from these surfaces by independently altering 363 the axes values of the surfaces in the interactive MATLAB plot window. The surfaces are 364 primarily color coded to indicate the gradient in solute solubility. Projections of these surfaces 365 manifest as grid lines (phase curves of ln(y)) on the axes planes. These plots indicate the major 366 and minute differences in the projected phase behavior put forth by the model equations. 367 Conveniently, even small or minute variations in a combinatorial pool of model equation 368 designs (derived from a parent equation) manifests acutely in the shape and color gradient of 369 the corresponding surface plots (Goos et al. 2011; Yamini and Moradi 2011; Cockrell et al. 370 2021). Further, literature (Schneider 1978; Mouahid et al. 2022) can be referred to make 371 accurate inferences regarding model specific phase behavior from these surfaces and 372 projections. However, a probable/possible approach (from the users' perspective) for gaining 373 satisfactory information from these surfaces (3D), its derivative plots and plane projections 374 (2D) is provided below. 375

Consider a set of model coefficient parameter estimates, from a DBSE model equation, 376 derived from empirical data from a (sufficiently) well modelled super/sub critical fluid 377 extraction process (for example, coffee or tea decaffeination) pertaining to a ternary system of 378 CO2/H2O solvent, Co-solvent (Ethanol or methanol) and solute (This ground truth data is 379 subject to availability and procurement by the user and is not provided here in/with this article). 380 Let this set of obtained estimates (which are highly process centric and equation specific) be 381 then used to plot the 3D surfaces and derivative plots (from this script). Naturally, due to the 382 process being sufficiently well modelled (as previously assumed), knowledge regarding the 383 projected Phase diagrams, vapor-liquid equilibrium behavior, maximum/optimal/desirable 384 solubility window and equilibrium points and planes is readily available, importantly reliable 385 and trustworthy for these estimates (ground truth), plots and the associated empirical data. Let 386 this information (again, not provided here with this article) be the ground truth and basis for 387 performing further comparative analysis using the MATLAB based plotting and visualization 388 script presented here in this article. Then the surfaces and 2D projections obtained by 389

implementing this visualization script for the same empirical data (and the model coefficient estimates) for a batch of DBSE model equations (existing/newly developed) can now be used to evaluate and glean information regarding the optimal window and other important associated attributes like the upper and lower critical end points, planes and edges associated with latency and the triple point. Further, vapor pressure curves and the data characteristic to the components (pure and mixture) in the ternary system can be identified and compared to this ground truth.

Generally, the qualitative and quantitative data regarding the latency, miscibility, 396 compression, crystallizability of the components in the reaction mix can be obtained from these 397 surfaces. Further, the identified solid-liquid-gas lines (by using cursor on the surface and 398 comparing point coordinates) describing boundaries of latency (or miscibility) projected on the 399 surface specific to the DBSE model equation(s) can also be compared to this (empirical) truth 400 and the error values quantify deviation and subtle / major differences. Similarly, values of slope 401 differentials (dP/dT, dT/dD and dP/dD) are easily computed from the surfaces for these 402 equations. The computed slope values could be used to identify upper and lower crossover 403 pressures bordering the retrograde solubility region in the phase diagrams for explaining / 404 utilizing retrograde solubility interference (Foster et al. 1991; Esmaeilzadeh and Goodarznia 405 2005; Kalikin et al. 2021). This is useful for screening newly designed DBSE model equations 406 for the maximum/optimal solubility window and the basis of which can further be used for 407 iteratively optimizing the optimal solubility window (by region specific selection), redesigning 408 customized, newer and efficient model equation alternatives. Overall, This Comparative 409 evaluation based on this ground truth is useful for selecting equations that project phase 410 behaviour with higher resolution and accuracy within the newly estimated/optimized optimal 411 solubility window. The Phase behavior projected from this particular newly selected equation 412 will now prove to be more beneficial for decision making and dynamic process optimization 413 (better than the present ground truth data). 414

Often, In Pilot / Production scale equipment, optimization is focused on Cost 415 effectiveness and dynamic feasibility factors (that influence cost effectiveness) including, but 416 not limited to, Process duration/Residence time, Resource consumption/availability, Climate 417 change/Ambient physical conditions, Hazard/Risk propensity among others. In such scenarios, 418 optimal solubility windows (based on current feasibility thresholds) can be estimated/optimized 419 (using this script) and implemented for large scale extraction. Conversely, Phase and VLE data 420 from the optimized surfaces derived from Model equations that are specifically designed (using 421 this work) for a particular process can potentially inform future decisions and process 422 trajectories for meeting optimal feasibility goals. Note that the maximum solubility window 423

(Not the optimal solubility window) depicted in Figure 4(b) is predicted and shown to lie around 424 the red regions (between 320K-340K and 30-32 MPa) by the tenth model equation (from the 425 same randomly mined sample of ten input equations). As pictorially shown, and explained, the 426 optimal solubility window is largely exploratory, process centric and will differ for a different 427 sample of equations for the same data (empirical ground truth) or other obvious differences in 428 physical parameters (residence time and quantity of reaction mix etc). To summarize, the plots 429 provide satisfactory, quantitative and qualitative knowledge regarding the phase behavior and 430 equilibria characteristic to the equations being studied, using this MATLAB based plotting and 431 visualization script. 432

433

434

#### Prediction of Solute Solubility: Machine Learning Algorithms

435

Multilayer Perceptron regression (MLP), K-Nearest Neighbours regression (KNN) and Support 436 Vector Regression (SVR) algorithms have been implemented using 'sklearn' package in python 437 in a single jupyter notebook. A toy data sample of 1000 randomly mined experiments are used 438 to illustrate the working of this jupyter notebook. The input parameters present in the toy data 439 sample are Temperature, Pressure and Density. The target / output / dependent variable is the 440 Solute solubility mole fraction. Additional parameters can be easily incorporated into the data 441 sample by simply concatenating them as columns after the Density data column in the input 442 data (Input\_Data.xlsx). The notebook initially provides description of the data using basic 443 statistical metrics (count, mean, standard deviation, minimum and maximum value), 444 Correlation values between the parameters and output, Heat map of correlation values and a 445 parameter pair plot for comparing the considered parameter pairs (combinations) on a chart. 446 These charts are shown in (Fig. 2 a - b). The results (graphs, errors and plots) and discussion 447 pertaining to each algorithm is provided in the subsequent paragraphs. 448



450

**Fig.** 5 Standard output from the jupyter notebook about the predictions and analysis of the Multilayer perceptron algorithm (a) Scatter Plot of Experimental (green) v/s Predicted (blue) values of solute solubility molefraction from the Multilayer perceptron algorithm. (b) Plot of residual values from the Multilayer perceptron algorithm. (c) Bar plot of error metrics of the predictions from the Multilayer perceptron algorithm.

Multilayer Perceptron regression (MLP) is the first algorithm implemented in this 456 Jupyter notebook. Data scaling (preprocessing) is performed using the 'MinMaxScaler' routine 457 before further transformation of the data. The data is then split (preprocessing) using the test-458 train-split routine. The results and the output obtained are illustrated in Fig. 5 a - c. Regression 459 model is built using the standard 'MLPregressor' routine. Hyperparameter optimization / tuning 460 is performed using the 'GridSearchCV' routine for the MLP algorithm. As explained, the space 461 for grid search for the hyperparameters (Number of hidden layers, activation functions, solvers, 462 learning rate) has to be defined in the beginning of the notebook for hyperparameter 463 optimization. Further, 5-fold cross validation is performed based on negated values of root 464 mean square error as the model scoring metric. The program performs tuning and the 465 hyperparameters of the best model are then used to refit and obtain the prediction output 466 (Schilling et al. 2015). Error metrics for this algorithm are (output) plotted (Fig. 5 c) separately. 467



469

468

**Fig.** 6 Standard output from the jupyter notebook about the predictions and analysis of the K-Nearest Neighbours algorithm. (a) Scatter Plot of Experimental (green) v/s Predicted (blue) values of solute solubility molefraction from the K- Nearest Neighbours algorithm. (b) Plot of residual values from the K- Nearest Neighbours algorithm. (c) Bar plot of error metrics from the K- Nearest Neighbours algorithm.

K-Nearest Neighbours regression (KNN) is implemented after MLP in the notebook.
As discussed, Data scaling was deemed unnecessary and has not been performed. However,
test train split is performed using the same routine as MLP. Further, The Hyperparameter K is
set to a random value of 3 for the toy data sample and can easily be changed / tuned based on

data and preference at the beginning of the notebook. Error metrics for the KNN algorithm is
plotted (Fig. 6 c) separately. Further insight regarding the model can be obtained from data,
hyperparameter optimization and previous literature (Soleimani Lashkenari and KhazaiePoul
2017).

483





Fig. 7 Standard output from the jupyter notebook about the predictions and analysis of the Support Vector Regression algorithm. (a) Scatter Plot of Experimental (green) v/s Predicted (blue) values of solute solubility molefraction from the Support Vector Regression algorithm.
(b) Plot of residual values from the Support Vector Regression algorithm. (c) Bar plot of error metrics from the Support Vector Regression algorithm.

Support Vector Machine Regression (SVR) algorithm is implemented next in the 490 notebook. Like before, data scaling is not performed so as to preserve pattern in the input 491 parameter space. Data has been split for model training using the test – train split routine like 492 before and can be easily adjusted. The choice of Kernel function hyperparameter is also tuned 493 using 'GridSearchCV' and the grid search space can be modified for this at the beginning of 494 the notebook. Five-fold cross validation is performed based on the negated root mean squared 495 error scoring metric and the kernel function associated with the best scoring model is then used 496 to refit and obtain the predictions (Tsirikoglou et al. 2017). Error metrics for SVR, like before, 497 is also plotted (Fig. 7 c) separately. 498





Fig. 8 Standard output from the jupyter notebook about the predictions and comparative 501 analysis of all three algorithms from each complete program run. (a) Combined plot of residual 502 values of all three machine learning algorithms (MLP, KNN, SVR) for comparison (b) Plot of 503 Prediction values of KNN and SVR algorithms with experimental solubility values. (c) Bar plot 504 of mean squared error values of all three machine learning algorithms (MLP, KNN, SVR) for 505 comparison. (d) Bar plot of root mean squared error values of all three machine learning 506 algorithms (MLP, KNN, SVR) for comparison. (e) Bar plot of percent absolute average relative 507 deviation values of all three machine learning algorithms (MLP, KNN, SVR) for comparison. 508 (f) Bar plot of mean absolute error values of all three machine learning algorithms (MLP, KNN, 509 SVR) for comparison. 510

Model fitting and prediction of all three algorithms (MLP, KNN and SVR) for the 511 sample data yielded results and the computed errors (MSE, RMSE, MAE and %AARD) are 512 plotted separately on bar plots (Fig. 8 c, d, e, f). The predictions v/s empirical data graph is also 513 plotted and exported by the notebook (Fig. 8 b). Likewise, residuals are also plotted for all three 514 algorithms (Fig. 8 a). Hyperparameter tuning for all three algorithms is implemented and other 515 intricate nuances (relative parameter importance, stacking of experiments based on specific 516 parameters/order) pertaining to the predictions can be easily made based on the best performing 517 algorithm and the input data (Feurer 2019). The parameter space, as previously explained, can 518 (only) be increased to explore and incorporate additional parameters like Residence time, 519 Mass/Volume of Raw Material, Viscosity of the Material, Melting point, Boiling Point, Total 520 polar surface area, Critical Temperature, Critical Pressure, Molecular Weight of solute, 521 percentage of co-solvent used, type of cosolvent (by scoring) etc. Therefore, Detailed 522 explanation regarding the obtained numerical output is unnecessary here since a toy data sample 523 with the standard (Temperature, Pressure and Density) parameters has been studied. The 524 notebook includes the best of the available plot commands and features (errors, functions, tables 525 etc) from standard python libraries for ease of use and assessment. The numerical data 526 predictions and analysis are also saved and exported to popular file formats (.xlsx). Importantly, 527

<sup>528</sup> users are cautioned against the usage of this notebook for actual experimental purposes as it can <sup>529</sup> be dangerous when used directly in a laboratory setting without proper consultation and <sup>530</sup> reasoning. The provided notebook is an efficient exploratory tool for data analysis and is very <sup>531</sup> useful for theoretical research, modeling (fitting), understanding and comparison. Overall, This <sup>532</sup> Jupyter notebook is a state-of-the-art predictive modeling and analysis tool using standard <sup>533</sup> Machine learning algorithms for obtaining prediction values of solute solubility mole fraction <sup>534</sup> from input parameter data.

535

536

#### Conclusions

537

This work describes software program scripts and presents their workflow as a comprehensive, 538 state of the art parameter estimation and predictive modeling tool for evaluating density based 539 semi empirical models (equations) and its associated data. Parameter estimation has been 540 implemented in a MATLAB based script using a customized version of the popular Ordinary 541 least squares estimation method. Further in this work, Visualization of phase behaviours 542 projected by preselected (sampled) model equations using a MATLAB based script has been 543 described. This visualization script produces three-dimensional surface plots in interactive 544 MATLAB windows based on the parameter estimates (computed from ordinary least squares 545 estimation). An approach for gleaning theoretical information regarding phase behaviour using 546 the surface plots is provided. Importantly, even subtle variations among model equations 547 acutely manifests in the shapes and color gradients of the projected surface plots and this makes 548 designing newer, robust, data specific/generalized equations easier. Standardized error and 549 scoring metrics have been computed at each appropriate stage in the workflow and is presented 550 in the form of plot illustrations. Importantly, the maximum solubility window is predicted to lie 551 somewhere around the red regions (probably between 320K-340K and 30-32 MPa) by the tenth 552 model equation (and is predicted for all the remaining input equations). The visualization 553 program is stand alone in that it fully functions when parameter estimates and parameter 554 boundaries for input equations are externally sourced. A Python based programming script is 555 also presented for predictive modeling of the associated input empirical data using three 556 Machine learning algorithms. Also, this notebook has been written to accommodate 'n' number 557 of other variables for improving the accuracy of the solute solubility predictions. This allows 558 users with diverse forms of data to easily make predictions, interpretations and reach 559 scientifically sound conclusions about the maximum/optimal solubility window. Further, user 560 defined hyperparameter tuning has been implemented for all three algorithms and has not been 561

entirely focused towards fitting the toy data sample (However, the presented error metrics are 562 desirably low). Therefore, it is strongly advised to use these program scripts for theoretical and 563 academic purposes since these scripts are under continuous development, refinement and 564 modification. The surfaces, plots and tables present in this article are the standard predictions 565 and analysis of outputs from these scripts based on a toy data and model equation(s) sample 566 (mined randomly from literature) and are not regarding any particular density based semi 567 empirical equation or published data. Hence, again, strong caution is advised against the usage 568 of any aspect of this work directly in an experimental setting without appropriate supervision 569 or reasoning. Importantly, a properly worded guide is provided for using this repository. Future 570 goals include deploying and testing this work on a GUI, established datasets, on temporal 571 variation, similar computational tools, and DBSE model equations. In summary, this work 572 postures a first of its kind, efficient computational tool in the form of program scripts for 573 evaluating/designing Density based semi empirical equations associated with super/sub critical 574 extraction process and data. 575

576

577 **Data Availability.** The Software programs are available in a GitHub Repository here 578 <u>https://github.com/Srinidhi-hub/DBSE-Evaluator.git</u>. The software programs are also 579 accessible upon request from the author.

580

581 Conflict of Interest. The Author declares that there are no conflicts of interests with any entity, 582 institution or individual regarding this study. The author also declares that this study/work did 583 not receive any funding from any source.

584

585 Acknowledgements. The Author is grateful to Prof. Dr. Rudiyanto Gunawan (Graduate 586 Academic/Research Advisor during my Graduate studies), Lab members, Professors, and 587 Colleagues at UB-CBE.

588

**Author Biography:** Srinidhi received his Master's degree in Chemical Engineering from The State University of New York at Buffalo (University at Buffalo) (CeDiD: 22G6-XJHD-SOI8). Earlier, he received his Bachelor's degree in Biotechnology Engineering from M. S. Ramaiah Institute of Technology. His Academic Interests include Mathematical modeling of Chemical and Biological Processes.



589				
590	Symbols			
591				
592	Т	Temperature	К	
593	Р	Pressure	MPa	
504	D	Density	$K\sigma/m^3$	
505	D	Density Desidual Sum of Squares	Kg/ III	
595	K	Residual Sum of Squares		
596	У	Solute Solubility Mole Fraction		
597				
598	Greek	z Letters		
599				
600	З	Error		
601	σ	Standard Deviation		
602	ρ	Density	$Kg/m^3$	
603				
604				
605		References		
005		interences		
606 607	Alwi	RS, Garlapati C (2021a) A new semi empirical model for the solub	ility of dyestuffs in	
608	S	upercritical carbon dioxide. Chemical Papers	75:2585–2595.	
609 610	h Alwi	https://doi.org/10.1007/s11696-020-01482-x RS_Garlanati C (2021b) A new model and estimation of thermodyna	mic parameters for	
611	the solubility of azobenzene and anthraquinone derivatives in supercritical carbon dioxide.			
612	Chemical Papers 75:4589–4610. https://doi.org/10.1007/s11696-021-01688-7			
613	Asgar	pour Khansary M, Amiri F, Hosseini A, et al (2015) Representing	solute solubility in	
614 615	a	nd Design 93:355–365. https://doi.org/10.1016/i.cherd.2014.05.004	gineering Kesearch	
616	Bartle	KD, Clifford AA, Jafar SA, Shilstone GF (1991) Solubilities of So	lids and Liquids of	
617	Ι	Low Volatility in Supercritical Carbon Dioxide. J Phys Chem Ref	Data 20:713–756.	
618	h D L'	https://doi.org/10.1063/1.555893	1 11 1 /	
619 620	Belits	er S V., Martens EP, Pestman WR, et al (2011) Measuring balance a	and model selection $20.1115 - 1129$	
621	ł	https://doi.org/10.1002/PDS.2188	ai 20.1115–112 <i>)</i> .	
622	Bian	XQ, Zhang Q, Du ZM, et al (2016) A five-parameter empirical mode	l for correlating the	
623	S	olubility of solid compounds in supercritical carbon dioxide. Fluid Pha	ase Equilib 411:74–	
624	8	30. https://doi.org/10.1016/j.fluid.2015.12.017	C 1 1 1	
625	Butlei	KI, Davies DW, Cartwright H, et al (2018) Machine learning	for molecular and $6-0.018-0.0337-2$	
627	Cock	ell C. Brazhkin V V., Trachenko K (2021) Transition in the supercrit	ical state of matter	
628	e	experimental evidence. Phys Rep. https://doi.org/10.1016/j.physrep.20	021.10.002	
629	Cunni	ngham P, Delany SJ (2022) k-Nearest Neighbour Classifiers - A Tuto	orial. ACM Comput	
630	S	Surv 54:1–25. https://doi.org/10.1145/3459665		

Esmaeilzadeh F, Goodarznia I (2005) Supercritical Extraction of Phenanthrene in the Crossover 633 Region. J Chem Eng Data 50:49–51. https://doi.org/10.1021/je049872x 634 Feurer M (2019) Automated Machine Learning. Springer International Publishing, Cham 635 Foster NR, Gurdial GS, Yun JSL, et al (1991) Significance of the crossover pressure in solid-636 fluid phase equilibria. Ind Eng Chem Res 30:1955-1964. supercritical 637 https://doi.org/10.1021/ie00056a044 638 Garlapati C, Madras G (2010) New empirical expressions to correlate solubilities of solids in 639 supercritical carbon dioxide. Thermochim Acta 500:123-127. 640 https://doi.org/10.1016/j.tca.2009.12.004 641 Goos E, Riedel U, Zhao L, Blum L (2011) Phase diagrams of CO2 and CO2–N2 gas mixtures 642 and their application in compression processes. Energy Procedia 4:3778–3785. 643 https://doi.org/10.1016/j.egypro.2011.02.312 644 Hawthorne SB (1990) ANALYTICAL-SCALE SUPERCRITICAL FLUID EXTRACTION. 645 Anal Chem 62:633A-642A. https://doi.org/10.1021/ac00210a722 646 Herrero M, Mendiola J, ... AC-J of C, 2010 undefined (2010) Supercritical fluid extraction: 647 applications. 1217:2495-2511. Recent advances and Elsevier 648 https://doi.org/10.1016/j.chroma.2009.12.019 649 Huang Z, Shi X, Jiang W (2012) Theoretical models for supercritical fluid extraction. J 650 Chromatogr A 1250:2-26. https://doi.org/10.1016/j.chroma.2012.04.032 651 Hunter JD (2007) Matplotlib: A 2D Graphics Environment. Comput Sci Eng 9:90–95. 652 https://doi.org/10.1109/MCSE.2007.55 653 Kalikin NN, Oparin RD, Kolesnikov AL, et al (2021) A crossover of the solid substances 654 solubility in supercritical fluids: What is it in fact? J Mol Liq 334:115997. 655 https://doi.org/10.1016/J.MOLLIQ.2021.115997 656 Khatib M El, de Jong WA (2020) ML4Chem: A Machine Learning Package for Chemistry and 657 Materials Science. arxiv.org. https://doi.org/https://doi.org/10.48550/arXiv.2003.13388 658 Knez Z, Skerget M, KnezHrnčič M (2013) Principles of supercritical fluid extraction and 659 applications in the food, beverage and nutraceutical industries. In: Separation, Extraction 660 and Concentration Processes in the Food, Beverage and Nutraceutical Industries. Elsevier, 661 pp 3–38 662 Kramer O (2013) K-Nearest Neighbors. pp 13–23 663 Lakshmi K, Mahaboob B, Rajaiah M, Narayana C (2021) Ordinary least squares estimation of 664 parameters of linear model. Journal of Mathematical and Computational Science 11:2015-665 2030. https://doi.org/10.28919/jmcs/5454 666 Matlab (1984) Matlab. The Mathworks.inc 667 Menke EJ (2020) Series of Jupyter Notebooks Using Python for an Analytical Chemistry 668 Course. J Chem Educ 97:3899–3903. https://doi.org/10.1021/acs.jchemed.9b01131 669 Mouahid A, Boivin P, Diaw S, Badens E (2022) Widom and extrema lines as criteria for 670 optimizing operating conditions in supercritical processes. J Supercrit Fluids 186:105587. 671 https://doi.org/10.1016/J.SUPFLU.2022.105587 672 Murtagh F (1991) Multilayer perceptrons for classification and regression. Neurocomputing 673 2:183-197. https://doi.org/10.1016/0925-2312(91)90023-5 674 Oliphant TE (2006) Guide to numpy, 2nd edn. Continuum Press 675 Pedregosa F, Michel V, Grisel O, et al (2011) Scikit-learn: Machine learning in Python. Journal 676 of Machine Learning Research 12:2825-2830. https://doi.org/http://scikit-677 learn.sourceforge.net/ 678

Dismuke, C R Lindrooth (2006) Ordinary least squares :Methods and Designs for Outcomes

631

632

Research

Rai A, Punase KD, Mohanty B, Bhargava R (2014) Evaluation of models for supercritical fluid 679 extraction. Int J Heat Mass Transf 72:274-287. 680 https://doi.org/10.1016/j.ijheatmasstransfer.2014.01.011 681 Reddy SN, Madras G (2011) A new semi-empirical model for correlating the solubilities of 682 solids in supercritical carbon dioxide with cosolvents. Fluid Phase Equilib 310:207–212. 683 https://doi.org/10.1016/j.fluid.2011.08.021 684 Roach L, Rignanese G, Fluids AE-... of S, 2023 undefined (2023) Applications of machine 685 Elsevier supercritical 202:896-8446. learning in fluids research. 686 https://doi.org/10.1016/j.supflu.2023.106051 687 Rovenski V (2010) Modeling of Curves and Surfaces with MATLAB®. Springer New York, 688 New York, NY 689 Schilling N, Wistuba M, Drumond L, Schmidt-Thieme L (2015) Hyperparameter optimization 690 with factorized multilayer perceptrons. Lecture Notes in Computer Science (including 691 subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 692 9285:87-103. https://doi.org/10.1007/978-3-319-23525-7\_6 693 Schneider GM (1978) Physicochemical Principles of Extraction with Supercritical Gases. 694 Angewandte Chemie International Edition in English 17:716-727. 695 https://doi.org/10.1002/anie.197807161 696 Schweidtmann AM, Esche E, Fischer A, et al (2021) Machine Learning in Chemical 697 Engineering: A Perspective. Chemie Ingenieur Technik 93:2029-2039. 698 https://doi.org/10.1002/cite.202100083 699 Selvaratnam B, Koodali RT (2021) Machine learning in experimental materials chemistry. 700 Catal Today 371:77-84. https://doi.org/10.1016/j.cattod.2020.07.074 701 Smola AJ, Schölkopf B (2004) A tutorial on support vector regression. Stat Comput 14:199– 702 222. https://doi.org/10.1023/B:STCO.0000035301.49549.88 703 Soleimani Lashkenari M, KhazaiePoul A (2017) Application of KNN and Semi-Empirical 704 Models for Prediction of Polycyclic Aromatic Hydrocarbons Solubility in Supercritical 705 Carbon Polycycl Aromat Compd 37:415-425. 706 Dioxide. https://doi.org/10.1080/10406638.2015.1129976 707 Tabernero A, del Valle EMM, Galán MÁ (2010) A comparison between semiempirical 708 equations to predict the solubility of pharmaceutical compounds in supercritical carbon 709 dioxide. Supercritical Fluids 52:161-174. Journal of 710https://doi.org/10.1016/j.supflu.2010.01.009 711 Tsirikoglou P, Abraham S, Contino F, et al (2017) A hyperparameters selection technique for 712 models. support vector regression Appl Soft Comput 61:139–148. 713 https://doi.org/10.1016/j.asoc.2017.07.017 714 W McKinney (2011) pandas: a foundational Python library for data analysis and statistics. 715 Python for high performance and scientific computing, . https://doi.org/http://pandas.sf.net 716 Yamini Y, Moradi M (2011) Measurement and correlation of antifungal drugs solubility in pure 717 supercritical CO2 using semiempirical models. Journal of Chemical Thermodynamics 718 43:1091–1096. https://doi.org/10.1016/j.jct.2011.02.020 719 720