1 **Predicting performance of in-situ microbial enhanced oil recovery process and**

2 **screening of suitable microbe-nutrient combination from limited experimental**

3 **data using physics informed machine learning approach**

4 Sree Pavan[a], K. Arvind[b], B. Nikhil[b], P. Sivasankar[a,*]

5 [a] Geo-Energy Modelling & Simulation Lab, Department of Petroleum Engineering,

6 Indian Institute of Petroleum & Energy, Visakhapatnam - 530003, India.

7 [b] Department of Mechanical, Chemical and Electronics Engineering,

8 OsloMet University, Oslo, Norway

9 [*] Corresponding author

10 Contact Details of Corresponding Author (P. Sivasankar)

11 Email ID: sivasankar.petro@iipe.ac.in

12 Ph: 91 9600043460

13

14

15

16

17 [Original Research Article]

18 Submitted to

19 Bioresource Technology

20 January 2022.

**Abstract**

To screen/identify suitable microbe, nutrient and reservoir for successful field implementation of in-situ MEOR technique, it is important to predict the oil recovery and quantify the relative importance of influencing parameters from limited experimental data. For this purpose, Physics-Informed Machine Learning (PIML) approach is adopted in this study, which is developed by integrating the physics-based and Machine Learning (ML) models. It is found that biosurfactant yield *w.r.t* nutrient ($Y_{PS}$), flow velocity and initial oil saturation ($S_{oi}$) are correspondingly the most influential microbial kinetic, operational and reservoir parameters. Higher oil recovery is achieved by selecting a microbe-nutrient-reservoir pair having higher $Y_{PS}$ and $S_{oi}$ values but with lower $Y_{XS}$ (microbial yield w.r.t nutrient) value. Among 12 ML models analysed, Neural network model had predicted the oil recovery relatively accurate ($R^2 \sim$ 0.98). Overall, this PIML approach helps to devise strategies for maximizing oil recovery at initial laboratory stage itself with limited experimental data.

Keywords: Microbial Enhanced Oil Recovery; Machine Learning; Biosurfactants; Modelling; Kinetics

**1. Introduction**

To meet the increase in global energy demand and to sustain crude oil production from depleting oil reservoirs, more than half of the crude oil that is left after primary and secondary recovery techniques must be recovered by suitable Enhanced Oil Recovery (EOR) techniques (Joshi et al., 2016). In relative to existing chemical EOR methods, In-situ Microbial Enhanced Oil Recovery (MEOR) method is an economical and environmental friendlier EOR technique (Joshi et al., 2016; Varjani and Upasani,

2

44     2016; Shibulal et al., 2018; Jeong et al., 2022). In in-situ MEOR process, exogeneous

45     (or) indigenous microbes are injected into the reservoir, which subsequently undergoes

46     metabolic activity within the reservoir by utilizing nutrients and producing bioproducts,

47     which consequently helps to recover the crude oil from the reservoirs (Joshi et al., 2016;

48     Varjani and Upasani, 2017; Shibulal et al., 2018; Markande et al., 2021). Though in-situ

49     MEOR technique inherits several advantages, it is not widely implemented in the field

50     across the globe as other chemical EOR techniques due to the existence of following

51     challenges (Nikolova and Gutierrez, 2020): (a) complexity in predicting the oil recovery

52     performance of in-situ MEOR technique; and (b) lack in quantifying the relative

53     importance of each influencing parameter on final oil recovery. Resolving these

54     challenges at initial lab investigation stage itself will correspondingly: help to decide

55     whether to implement in-situ MEOR technique in the given reservoir or not and to

56     identify/screen the suitable microbe-nutrient-reservoir combination for attaining better

57     oil recovery; and assist in development of strategies for optimizing the oil recovery.

58        To evaluate the oil recovery performance of in-situ MEOR process, earlier,

59     several core flooding experimental studies (Joshi et al., 2016; Varjani and Upasani,

60     2016; Shibulal et al., 2018) and physics based computational modelling studies (Nielsen

61     et al., 2016; Sivasankar and Kumar, 2016, 2017, 2019; Jeong et al., 2021, 2022) were

62     performed.  However, performing core flooding experimental studies to identify/screen

63     a suitable microbe-nutrient-reservoir combination from several available combinations

64     makes experimental approach an expensive and time-consuming exercise. Though

65     physics-based models can provide better prediction of oil recovery with physically

66     consistent results, but it is computationally intensive to perform uncertainty

67     quantification and optimization studies as it requires to solve the non-linear equations

68  for several simulation runs (Thanh et al., 2020; Karniadakis et al., 2021). Moreover, it is

69  also unfeasible to quantify the relative importance of each influencing parameter on oil

70  recovery by both experimental and physics-based modelling approach as it requires

71  several experiments. Recently, Machine Learning (ML) models/algorithms are

72  increasingly used for different bioprocess applications to predict and optimize its

73  performance (Cruz et al., 2012; Tang et al., 2021; Zhang et al., 2021; Wang et al.,

74  2022). With the availability of large input and output datasets, ML models can quickly

75  predict the outcome of complex problems and quantify the relative importance of each

76  input parameters, which is otherwise difficult by using only physics-based models

77  (Thanh et al., 2021; Tang et al., 2021). However, with the limited availability of data

78  from experimental and field studies, it will not be feasible to apply ML

79  models/algorithms alone as it may predict physically inconsistent results with lesser

80  accuracy. Hence the requirement to have a quick and physically consistent results from

81  limited observed/experimental data with better accuracy is achieved by integrating both

82  the physics informed model and data driven ML model into a single hybrid model

83  called Physics Informed Machine Learning (PIML) model (Thanh et al., 2020;

84  Karniadakis et al., 2021). In recent times, PIML modelling approach is gaining

85  popularity because of its ability to accommodate the merits of both physics-based model

86  and ML model in a single model, while mitigating their respective drawbacks. Recently,

87  PIML modelling approach have been successfully used for different applications (Thanh

88  et al., 2020; Karniadakis et al., 2021; Liu et al., 2021). However, the use of PIML

89  approach for in-situ MEOR application have not been explored yet at least to the

90  authors knowledge.

SOMETHING ABOUT PIML WRITE ABOUT EXPLAINABLITY,

INTERPRETABLE AND PHYSICALLY CONSISTENT…

Hence the novelty of the present work is in introducing the PIML approach for in-situ

MEOR application to predict its oil recovery performance and to quantify the relative

importance of parameters influencing the oil recovery using limited experimental data.

In particular, the objectives of the present work are: (a) to develop a framework to

integrate the physics based model and ML model into a single PIML model for

generating a large, relevant and physically consistent data sets from limited

experimental data; (b) to quantify the relative importance of each parameter on

influencing the final oil recovery using PIML approach, and subsequently to identify the

critical kinetic and operational parameters influencing the oil recovery; and (c) to

identify the suitable ML model among 12 different ML models that shall be used

directly in PIML approach for predicting the oil recovery performance.

The present PIML approach study will help the end-users: to quickly select a

favourable microbial-nutrient-reservoir combination from several other available

options; to decide whether to implement in-situ MEOR technique in a particular

reservoir or not; and to devise operational strategies for maximizing the oil recovery.

**2. Materials and Methods**

In the present study, PIML approach is developed by combining the physics-

based model and ML model into a single model. Initially, laboratory experiments are

performed to determine the microbial kinetic and reservoir properties data. Based on

114 these limited experimental data, physics-based models for microbial kinetic and oil

115 recovery processes are developed. From these physics-based models, large, physically

116 relevant input and output data sets are generated. Using these large datasets, the ML

117 models are then trained and tested to quantify the relative importance of each parameter

118 and to predict the oil recovery quickly. The methodology for developing this PIML

119 approach is presented in detail in this section and briefed in Figure 1.

120 [Figure 1]

121 ### 2.1 Classification and collection of input parametric data

122 In the present study, 13 input parameters are considered. The corresponding

123 values of these input parameters constitutes the input parametric data. In the present

124 study, the input parametric data are classified as: (i) microbial kinetic parametric data,

125 (ii) operational parametric data, and (iii) reservoir parametric data, based on the

126 corresponding properties of microbes, nutrients, operational and reservoir conditions.

127 *2.1.1 Collection of input microbial kinetic parametric data from experimental studies*

128 In the present study, the microbial kinetic parameters that are considered as

129 input are maximum microbial growth rate [$U_{max}, (h^{-1})$], yield of microbes *w.r.t* sucrose

130 ($Y_{XS}$), yield of biosurfactants *w.r.t* sucrose ($Y_{PS}$) and Monod half saturation coefficient

131 ($K_{XS}, (gl^{-1})$]. The corresponding values of these microbial kinetic parameters are

132 considered as input microbial kinetic parametric data. In the present study, the input

133 data for all these microbial kinetic parameters are sourced from the experimental studies

134 of Sivasankar et. al., 2016, in which, *Pseudomonas putida* MTCC 2467 was used as

135 microbe, while sucrose and ammonium sulphate were used as carbon and nitrogen

136 source nutrient, respectively. In that study, at pH 8 condition, experiments on microbial

137　　growth, nutrient utilization and biosurfactant production were carried out to determine

138　　the values of microbial kinetic parameters for predicting the oil recovery.

139　　*2.1.2 Collection of input operational and reservoir parametric data*

140　　　　In the present study, the operational parameters that are considered as input are

141　　mean flow velocity of injection fluid within reservoir $[u, (mh^{-1})]$, viscosity of injection

142　　fluid $[\mu_w, (Nhm^{-2})]$, initial/injection concentration of microbes $[X_i, (gl^{-1})]$,

143　　initial/injection concentration of sucrose $[S_i, (gl^{-1})]$, initial/injection concentration of

144　　ammonium sulphate $[A_i, (gl^{-1}),]$ and resident time $[T_r, (h)]$.  These input operational

145　　parameters are controlled by the operators/scientists in the field/laboratory during the

146　　implementation of in-situ MEOR technique. Finally, the reservoir fluid-rock parameters

147　　that are considered as input parameters in the present study are initial residual oil

148　　saturation $[S_{ori}, (fraction)]$, irreducible water saturation $[S_{wir}, (fraction)]$ and initial or

149　　maximum oil-water Interfacial Tension (IFT) at the start of EOR $[\sigma_{max}, (mNm^{-1})]$. In

150　　the present study, the input data for all these operational and reservoir rock-fluid

151　　parameters are sourced from Sivasankar et al., 2016. Table 1 presents the sourced data

152　　or reference value of all these input parameters. It is to be noted that for each input

153　　parameter, only one reference value is available either from experiments or other

154　　sources, which will be insufficient for applying the ML algorithms.

155　　　　　　　　　　**[Table 1]**

156　　**2.2 Generation of large input and output datasets from physics-based model**

157　　　　In the present study, percent of oil recovery is the only parameter considered as

158　　output parameter. This output oil recovery parameter is influenced by all the input

159　　parameters (Sivasankar et al., 2016) that are mentioned in Tab. 1. In order to apply

Machine Learning (ML) algorithms for predicting the output oil recovery, large data

sets of input and output parameters are required to train and test the ML algorithms.

However, the availability of input and output parametric data from laboratory

experiments and other sources are limited (as presented in Tab. 1), which is inadequate

to implement the ML algorithms. Hence in the present study, large datasets of input and

output parameters are generated synthetically (data augmentation) for training and

testing the ML algorithms. Data augmentation is a mathematical method to synthesise

more data from the known (experimental) data when there is data insufficiency. The

methodology adopted in the present study for generation of input and output data is

similar to the method earlier adopted by Thanh et al., 2020, and it is outlined in sec.

2.2.1. and sec. 2.2.2.

*2.2.1 Generation of large input datasets from sourced reference values*

The reference value of input microbial kinetic parameters that are presented in

Tab. 1 are specific only to a particular temperature, pH, salinity, and pressure conditions

at which experiments were conducted. However, in actual reservoir fields, the reference

value of input parameters mentioned in Tab. 1 varies significantly due to the existence

of heterogeneity, resulting in uncertainty (Ansah et al., 2020; Thanh et al., 2020). Hence

accounting for this uncertainty, and to make the present model to be applicable for

wider variations in input parametric data during its field implementation, a 50%

Standard Deviation (SD) is considered to all the input parameter values (Thanh et al.,

2020). The resultant value range for each of these input parameters after considering the

SD is presented in Tab. 1. Subsequently, large datasets of about 10000 values (i.e., data)

for each of the input parameter is generated between their corresponding value range by

dividing it in equal intervals following the uniform distribution.

184    *2.2.2 Generation of large output data sets using physics-based model*

185    The output data on oil recovery is dependent on all the input parametric data.

186    Hence, to generate the large data sets of output parameter (oil recovery, %) for training

187    and testing of ML algorithms, the physics-based model (Eqs. 1- 8) (Sivasankar and

188    Suresh Kumar, 2016) which is dependent on all input parameters is simulated several

189    times using the generated input datasets. In the present physics-based model, the

190    microbial kinetic model (Eqs. 1 – 4) simulates: growth kinetics of microbes (Eq. 1);

191    nutrient utilization kinetics (Eq. 2); biosurfactant production kinetics (Eq. 3); and

192    Monod's kinetics (Eq. 4). While the oil recovery model (Eqs. 5 – 8) simulates: IFT

193    reduction by produced biosurfactants (Eq. 5); increase in Capillary Number due to IFT

194    reduction (Eq. 6); decrease in oil saturation due to increase in Capillary Number (Eq. 7);

195    and the final percent of oil recovery (Eq. 8), which is the output and target data. Based

196    on this obtained oil recovery data, performance evaluation of MEOR technique and

197    screening of suitable microbe-nutrient-reservoir combination are carried out.

198    $dX/dt = \mu_x.X$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀(1)

199    $dS/dt = -\mu_x.X/Y_{XS};$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀$dA/dt = -\mu_x.X/Y_{XA}$ ⠀⠀⠀⠀(2)

200    $dP/dt = (Y_{PS}/Y_{XS}).\mu_x.X$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀(3)

201    $\mu_x = U_{max}.\{(S/K_{XS} + S) + (A/K_{XA} + A)\}$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀(4)

202    $\log(\sigma^*) = \log(\sigma_{min}) + \log(\sigma_{max}/\sigma_{min}).\{(P - P_{max})/(P_{max} - P_{min})\}$ ⠀⠀⠀(5)

203    $N_{ca}^* = u_w\mu_w/\sigma^*$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀(6)

204    $S_o = \left[\frac{-\tanh(v_1(N_{ca}^*)-v_3)+1+v_2}{-\tanh(v_1(N_{ca}^0)-v_3)+1+v_2}\right]S_o'$ ⠀⠀⠀⠀⠀⠀⠀⠀$S_w = 1 - S_o$ ⠀⠀⠀(7)

9

205    $Oil\ recovery, \% = \{(S_w - S_{wir})/(1 - S_{wir})\} \times 100$ $\hspace{2cm}$ (8)

206    $\hspace{1cm}$ In Eqs. (1 - 8), the terms: *X, S, A* and *P* represents the concentration of microbes,

207    sucrose ammonium sulphate and produced biosurfactant, respectively in $gl^{-1}$; $\mu_x$

208    represents the microbial growth rate in $h^{-1}$; $K_{XA}$ represents the half-saturation constant

209    of ammonium sulphate in $gl^{-1}$; $Y_{XA}$ represents the yield of microbes *w.r.t* ammonium

210    and sulphate, r; $Y_{PS}$ represents the yield of biosurfactant *w.r.t* sucrose; $N_{ca}$ represents

211    the updated IFT ($mNm^{-1}$) and Capillary Number, respectively; $P_{min}$ and $P_{max}$

212    represents the minimum and maximum biosurfactant concentration, respectively in

213    $gl^{-1}$; $\sigma^*$ and $\sigma_{min}$ represents the updated IFT and minimum IFT, respectively in

214    ($mNm^{-1}$); $S_o$ and $S_w$ represents the saturation of oil and water, respectively in fraction;

215    and $\nu_1, \nu_{2,} \nu_3$ are the constants.

216    $\hspace{1cm}$ By performing one simulation job of physics-based model from Eqs. 1 - 8, one

217    output data on oil recovery is generated. Hence, in the present study, to generate a large

218    database of output data, ten thousand (10000) simulation jobs were performed which

219    resulted in generation of 10000 output data on % oil recovery. While, in each simulation

220    job, the input value (data) of different input parameters that are required are sampled

221    randomly from the generated input database using Latin Hyper-Cube Sampling (LHS)

222    technique (Thanh et al., 2020). In some simulation jobs, the set of input data have not

223    generated a valid positive output data (i.e., % of oil recovery), and such data are

224    excluded from the training and testing of ML algorithms. Figure 2 shows the frequency

225    distribution of all the input and output data values that were considered in the present

226    study for training and testing of different ML algorithms.

227    $\hspace{3cm}$ [Figure 2]

228    *2.3 Machine Learning Models*

229    Subsequent to the generation of large sets of input and output data, the

230    interaction strength (or) sensitivity of all 13 input parameters on the output oil recovery

231    is quantified by using Pearson Correlation Coefficient (PCC) and Spearman Rank

232    Corelation Coefficient (SRCC) values. PCC value measures the linear relationship

233    between two variables and SRCC value measures the monotonic relationship between

234    two parameters. Both PCC and SRCC values range from -1 to 1. Positive correlation

235    value between two parameters indicates that parameters are directly proportional, and

236    *vice versa*. Magnitude of the correlation indicates the strength of the relationship

237    between the two parameters. Higher the magnitude of correlation coefficient value,

238    higher is the association strength between the two parameters. Determination of PCC

239    and SRCC values helps to quantify the influence of different input parameters on output

240    oil recovery, which shall be used to screen the suitable microbes, nutrients and

241    reservoirs at the laboratory experimental stage for MEOR field implementation.

242    *2.4 Prediction of relative importance score to quantify the significance of input*

243    *microbial kinetic, operational and reservoir parameters on output oil recovery*

244    In the present work, feature importance study is carried out to quantify the

245    relative importance of each input parameter on the output oil recovery using Random

246    Forest Classifier ML algorithm (Keprate and Ratnayake, 2019) in the present PIML

247    framework. This ML algorithm has been trained and tested using the input and output

248    parameter datasets that are generated from physics-based model (as described in sec

249    2.3). This feature importance study computes the Relative Importance (RI) score for

250    each input parameters in fraction, where its summation will be 1. Hence, RI score of an

251 input parameter quantifies the significance (or) importance of that input parameter on

252 influencing the output oil recovery in relative to other input parameters. Determination

253 of this RI score for all input parameters will helps to identify the critical input

254 parameters influencing the output oil recovery, which subsequently guide the future

255 operation. In the present work, the results from the feature importance study would

256 helps: (a) to identify the input parameters that are most and least important for

257 predicting the oil recovery, which subsequently helps to identify the input parameters

258 which exhibits higher and lower influence on the output oil recovery; (b) to identify the

259 critical input parameters that shall be optimized for improving the efficiency of oil

260 recovery; and (c) to determine the weightage functions of all input parameters, which

261 shall be used to screen the suitability of MEOR technique among other EOR techniques

262 and to identify the right combination of microbe-nutrient pair for attaining better oil

263 recovery during its field implementation.

264 *2.5 Prediction and evaluation of different machine learning algorithms for MEOR*

265 *applications from lab data*

266 In the present study, Machine Learning (ML) model which is integrated within

267 the PIML approach is used to predict the output oil recovery. CRISP-DM methodology

268 was used for performing data mining and predicting the output oil recovery from input

269 parameters (Keprate and Ratnayake, 2019). The large data sets of input parameter data

270 and output data that are required for training and testing the ML model are sourced from

271 physics-based model which is embedded within the PIML approach (the procedure for

272 data generation using physics-based model is presented in sec. 2.2). As there are

273 different ML models available to do the prediction, it is necessary to identify the most

274 accurate and suitable ML model that can be used in the PIML approach by the end-users

275 (researchers/scientists in the laboratory) for predicting the oil recovery. Hence, in the

276 present study, 12 different ML models/algorithms are used in the PIML approach to

277 determine its accuracy in predicting the output oil recovery. The 12 different ML

278 models/algorithms that were used in the present study are K-Nearest Neighbours

279 (KNN), Decision Trees, Lasso, Ridge, Linear Regression, Random Forests, ADA Boost

280 Regression, Gradient Boosting, Gaussian Process Regression, Polynomial Regression,

281 Support Vector Regression (SVR) and neural networks.

282     For all these 12 ML models adopted in the PIML approach, the input parameter

283 data was normalised, and was subsequently split into training data sets and test data sets

284 in the ratio of 7:3 for training and testing of the ML model used. k-fold cross validation

285 was performed on the training set by setting k = 10, and the best model is then evaluated

286 on the test data set. In the present work, all the 12 ML models were trained using

287 training data sets, and its prediction performance were compared based on 3 metrics,

288 namely, Root Mean Square Error {$RMSE$; eq. (9)}, Coefficient of Determination {$R^2$;

289 eq. (10)} and Explained Variance Score {$EVS$; eq. (11)}.

290 $$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}} \tag{9}$$

291 $$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \tag{10}$$

292 $$EVS = 1 - \frac{Var(y_i - \hat{y}_i)}{Var(y_i)} \tag{11}$$

293 In the Eqs. (9 - 11), where:$y_i$ represents the actual % oil recovery determined from

294 physics-based model; $\hat{y}_i$ represents the predicted value of % oil recovery determined

295 from ML model; $\bar{y}$ represents the mean value of $y_i$; $n$ represents the number of

296 samples; and *Var* represents the variance. *RMSE* is a measure of accuracy, and lower

297 values indicate better fit of data. $R^2$ and $EVS$ measures proportion to which a

298 mathematical model accounts for variation of a given data set. The ML model having

299 values of $R^2$ and $EVS$ closer to 1 is the most accurate and suitable model that shall be

300 used to predict the % oil recovery for in-situ MEOR application.

301 **3. Results and Discussion**

302 *3.1 Validation of physics based microbial kinetic model*

303       The validity of the physics-based microbial kinetic model that is used in the

304 present study is verified by comparing the present numerical model results with the

305 experimental data. From Fig. 3(a - c), it is observed that the present model results

306 (microbial, nutrient and bio-surfactant concentrations *w.r.t* time) is in good agreement

307 with the experimental data. As the present adopted model is validated, it is subsequently

308 used to generate large datasets.

309                             [Figure 3]

310 *3.2 Quantifying the influence of input parameters on output oil recovery*

311                             [Figure 4]

312       Figure 4a shows the PCC and SPCC values in a matrix form that represents the

313 interaction (or) association strength between any two parameters involved in the MEOR

314 process. Fig. 4b specifically presents the PCC and SPCC values (i.e., interaction

315 strength) between all the input parameters with the output oil recovery parameter.

316 Results from Fig. 4a and Fig. 4b reveals that the input parameters, $Y_{PS}, u, S_{ori}, \mu_w,\ X_i,$

317 $A_i, S_i, U_{max},\ S_{wir}\ and\ T_r$ are directly proportional to the oil recovery, while the input

318  parameters $K_{XS}, Y_{XS}$ and initial IFT are inversely proportional to the oil recovery.

319  These results are consistent with the reality, which validates the results shown in Fig. 4.

320  It is observed that among all these input parameters, the parameter, $Y_{PS}$ had

321  strongly associated with the output oil recovery while compared to all the other input

322  parameters. This illustrates that the yield value of biosurfactants *w.r.t* sucrose ($Y_{PS}$) is

323  the dominant parameter that significantly influences the output oil recovery. Moreover,

324  it is also inferred that the oil recovery performance of MEOR process increases with

325  increase in $Y_{PS}$ value, which means that with the increase in utilization of nutrients for

326  biosurfactant production, the oil recovery increases. This obtained result corroborates

327  with the earlier results of Sivasankar and Suresh Kumar, 2019, in which, it is reported

328  that $Y_{PS}$ parameter significantly influences the oil recovery compared to other kinetic

329  parameters. From Fig. 4a and Fig. 4b, it is also observed that among the negatively

330  correlated input parameters (*i.e.*, parameters that are inversely proportional to the oil

331  recovery), $Y_{XS}$ is the input parameter that is strongly associated with the output oil

332  recovery. This illustrates that lower the value of $Y_{XS}$ (*i.e.*, less nutrient is utilized for the

333  growth of microbes), higher is the oil recovery. Hence, the study reveals that the higher

334  oil recovery is attained by selecting a microbe-nutrient pair that have higher value of

335  $Y_{PS}$ and lower value of $Y_{XS}$. Based on these observations made on $Y_{PS}$ and $Y_{XS}$ values, it

336  shall be finally correlated that the ratio between $Y_{PS}$ and $Y_{XS}$ (*i.e.*, $Y_{PS}/Y_{XS}$) values for a

337  microbe-nutrient pair needs to be higher to achieve better oil recovery. Thus, based on

338  the determination of PCC and SPCC values, it is concluded that: (a) $Y_{PS}$ and $Y_{XS}$ are the

339  two input parameters that significantly influences the output oil recovery; and (b) oil

340  recovery could be maximized by selecting a microbe-nutrient pair having higher

341  $Y_{PS}/Y_{XS}$ value at initial laboratory investigation stage itself.

342    *3.3 Application of PIML approach for identifying critical input microbial kinetic,*

343    *operational and reservoir parameters influencing the output oil recovery*

344                            [Figure 5]

345        Figure 5a presents the Relative Importance (RI) score or relative strength of all

346    the 13 input parameters on influencing the output oil recovery. The RI score was

347    determined by performing feature importance study. As all the 13 input parameters

348    involved in the feature importance study were selected (or) sourced from the physics-

349    based model (Eqs. 1- 8), hence, all these input parameters are relevant and have some

350    influence on deciding the output oil recovery. This is evident from Fig. 5a, which shows

351    that each of the 13 input parameters have a non-zero RI score. Thus, in the present

352    PIML approach, all the 13 input parameters are considered for the training, testing, and

353    implementation of all ML algorithms (models) for predicting the output oil recovery.

354        It is observed from Fig. 5a that among all the input parameters, $Y_{PS}$ has the

355    highest RI score of 0.168, hence, it is the most critical input which significantly

356    influences the output oil recovery. In order to exclusively understand the relative

357    importance of microbial, operational and reservoir parameters on deciding the output oil

358    recovery, correspondingly, Figs. 5b, 5c and 5d are plotted. It is understood from Fig. 5b

359    that among the input parameters that are related to microbes and nutrients (i.e., $Y_{PS}, Y_{XS},$

360    $S_i, A_i, U_{max}, X_i, K_{XS}$), $Y_{PS}$ and $Y_{XS}$ are relatively the most influential input parameters

361    with RI score of 0.168 and 0.1, respectively. While, $K_{XS}$ is relatively the less significant

362    input kinetic parameter on deciding the percent of output oil recovery with RI score of

363    0.014. It is also observed from Fig. 5b that compared to injection concentration of

364    microbes, the injection concentration of nutrients (both, carbon and nitrogen source)

365    into the reservoir has relatively higher impact on deciding the output oil recovery. This

366    implies that for maximizing the oil recovery, continuous supply of nutrients to the

367    microbes need to be ensured for microbes to undergo metabolic activity within the

368    reservoir (*i.e.*, to produce biosurfactants) and recover the oil.

369         Fig. 5c shows the relative importance of operational parameters ($u, \mu_w, T_r$) on

370    influencing the output oil recovery. It is observed from Fig. 5c that though all

371    operational parameters influence the output oil recovery, fluid velocity ($u$) is the input

372    operational parameter that influences the output oil recovery relatively more, and

373    closely followed by the viscosity of injection water ($\mu_w$) parameter. This obtained

374    results are in accordance with the physics-based concept of Capillary Number, in which,

375    the viscous force (*i.e*, product of $u$ and $\mu_w$) must be higher for achieving higher oil

376    recovery. Thus, the results from Fig. 5c implies that more oil could be recovered during

377    field implementation of in-situ MEOR technique by optimizing the injection velocity of

378    the microbial slug (*i.e*, mixture of microbes, nutrients and water) and by increasing the

379    water viscosity using biopolymer producing microbes during in-situ MEOR application.

380         Figure 5d presents the relative importance scores of different parameters (*i.e.*,

381    $S_{ori}, S_{wir}, \sigma_i$) related to the fluids present within the reservoir. By correlating the results

382    from Fig. 5d and from Fig. 4b, it is inferred that among the fluid parameters, the initial

383    residual oil saturation parameter ($S_{ori}$) is the most significant parameter influencing the

384    oil recovery, and the oil recovery will be higher in reservoirs that has higher value of

385    $S_{ori}$. This finding is in good agreement with the earlier physics-based simulation studies

386    (Sivasankar et al., 2016) which states that the oil recovery performance increases with

387    the increase in initial residual oil saturation. Hence, based on this finding from Fig. 5d it

388    is suggested that the oil recovery performance of MOER technique could be improved if

389    the MEOR technique is implemented at an earlier stage of oil production (i.e., along

390    with secondary recovery stage), during which the oil saturation will be relatively higher

391    compared to the later stage (i.e., at tertiary recovery stage).

392         As the results presented in Fig. 5 (a – d) validates with the physics-based model

393    results, it has been affirmed that the physics has been infused into the ML model in the

394    present PIML approach. Hence the results obtained from this PIML approach can be

395    used to draw physical insights, based on which, suitable strategies can be developed for

396    maximizing the oil recovery.

### 3.4 Application of PIML approach to screen suitability of in-situ MEOR technique and to identify suitable microbe-nutrient for in-situ MEOR implementation

397    *3.4 Application of PIML approach to screen suitability of in-situ MEOR technique*

398    *and to identify suitable microbe-nutrient for in-situ MEOR implementation*

399         The RI score of each input parameter presented in Fig. 5a also correspondingly

400    represents the weightage factor of each input parameter. Based on this weightage factor,

401    the selection score of in-situ MEOR technique is calculated. This selection score helps

402    in initial screening of in-situ MEOR technique among other EOR techniques for field

403    implementation. The EOR technique that possess the highest selection score will be

404    considered further for field implementation. Earlier, the selection sore for in-situ MEOR

405    technique was calculated based on the reservoir fluid and rock properties, and neglected

406    the consideration of microbial kinetic and operational parameters, which may mislead

407    the entire EOR screening process for field implementation. However, with the PIML

408    approach adopted in the present work, the selection score for in-situ MEOR technique is

409    calculated by including both microbial kinetic parameters ($Y_{PS}, Y_{XS}, K_{XS}$) and

410    operational parameters along with reservoir fluid and rock parameters. Thus, the present

411    work would enhance the accuracy in screening of in-situ MEOR technique, which

18

412     subsequently would help the end users to make better decision on selecting a suitable

413     EOR for field implementation. The selection score for an EOR technique is calculated

414     by using the formula, $Selection\ Score = \sum_{i=1}^{n} w_i a_i$. Here, $i$ represents the input

415     parameters; $n$ represents the total number of input parameters; $a_i$ represents the

416     accuracy factor of input parameter, $i$, and its value varies between 0 and 1. Accuracy

417     factor value represents the measure of closeness of that input parameter value with the

418     most favourable value range; $w_i$ represents the weightage function of the input

419     parameter, $i$, and its value varies between 0 and 1. This weightage factor represents the

420     relative importance of that input parameter influencing the output parameter.

421        In calculation of selection score for present in-situ MEOR technique, the

422     weightage factor, $w_i$, of different input parameters, $i$, are same as the RI score of

423     different input parameters as shown in Fig. 5a. Hence, based on the determined

424     weightage factor (*i.e*., RI score) for all the 13 input parameters, the selection score for

425     the in-situ MEOR technique shall be calculated by using Eq. (12).

426     $Selection\ Score = 0.168\ a_{Yps} + 0.146\ a_{Sori} + 0.14\ a_u + 0.125\ a_{\mu w} + 0.114\ a_{Ai} +$

427                  $0.1\ a_{Yxs} + 0.08\ a_{Si} + 0.047\ a_{Umax} + 0.02\ a_{Tr} + 0.018\ a_{Swir} +$

428                  $0.016\ a_{Xi} + 0.014\ a_{Kxs} + 0.012\ a_{\sigma i}$            (12)

429     The value of accuracy factor values of each input parameter ($a_i$) are determined from

430     lab experiments. The value of $a_i$ varies case-to-case basis, and its value depends on the

431     nature of reservoir and microbe-nutrient pair used and the operational conditions

432     adopted. Upon calculation of $a_i$ from initial experiments, the selection score for in-situ

433     MEOR technique shall be quickly calculated using Eq. (12), which will subsequently

434     help to screen the suitability of in-situ MEOR technique among other EOR techniques

435     for field implementation at the initial laboratory investigation itself. In addition to it, the

436     selection score presented in Eq. (12) also helps to identify the suitable microbe-nutrient

437     combination among several available combinations for attaining better oil recovery at

438     the initial laboratory investigation itself. The microbe-nutrient combination that have

439     highest selection score value will recover relatively more oil from the reservoir for a

440     given reservoir and operational conditions. Thus, it is concluded that the RI score

441     determined from feature importance study in present PIML approach will: (a) help to

442     screen the suitability of in-situ MEOR technique among other EOR techniques for field

443     implementation; and also (b) helps to screen the suitable microbe-nutrient combination

444     for successful implementation of in-situ MEOR technique in the field at the initial

445     laboratory investigation itself.

446     ***3.5 Application of different ML algorithms in the PIML approach to predict the oil***

447     ***recovery performance of in-situ MEOR technique***

448                                 [Figure 6]

449        Figure 6 shows the oil recovery (in %) predicted by different ML models used in

450     the PIML approach against the benchmark (actual) results which are obtained from

451     physics-based models. The most accurate ML model with better prediction capability

452     will have the scatter plot points lying closer to the line equation $y' = x$ (here, $y'$ and $x$

453     are benchmark and predicted values, respectively), and correspondingly will have $R^2$

454     and *RMSE* value closer to 1 and 0, respectively. While, for the ML model with least

455     accuracy, the scatter plot points are spread widely from the line equation $y' = x$, and it

456     will also have relatively lower $R^2$ and relatively higher *RMSE* value. The $R^2$, *RMSE* and

457     *EVS* values of all the 12 ML algorithms that were used in the present PIML approach

458     study is presented in Table 2.

459       [Table 2]

460       Based on the results presented in Fig. 6 and Tab. 2, it is inferred that among all

461       the 12 ML algorithms/models that are used in the present PIML approach, Neural

462       Networks ML model had performed better in predicting the output oil recovery ($R^2 =$

463       0.9873, $RMSE = 0.7145$). The neural networks ML algorithm/model outperforms other

464       ML algorithms in prediction because it can implicitly detect complex non-linear

465       relationships between dependent and independent variables, and it also have the ability

466       to detect all possible interactions between the input variables. Followed by the neural

467       network model, it is observed that the Support Vector Regression (SVR) is the second-

468       best ML model that can better predict the oil recovery ($R^2 = 0.9644$; $RMSE = 1.184$).

469       The main advantage of SVR model is that it is less susceptible to outliers than other

470       data-driven models but it's harder to manually tune hyperparameters. Next to SVR

471       algorithm, it is found that the 4th degree Polynomial Regression model had predicted the

472       oil recovery better ($R^2 = 0.963$; $RMSE = 1.26$) as it has the ability to better map the

473       non-linear relationship between the input and output variables. Amongst all the 12 ML

474       models that were used in the present PIML approach for oil recovery prediction, it is

475       found that K-Nearest Neighbours ML model is the least accurate model ($R^2 = 0.3698$;

476       $RMSE = 3.369$). Thus, from the present study, it is concluded that to predict the oil

477       recovery performance of in-situ MEOR technique at initial lab stage, the Neural

478       Network is the best ML algorithm that need to be used in the PIML approach.

479       *3.5 Case study on application of PIML approach for screening of suitable microbe-*

480       *nutrient combination for in-situ MEOR implementation*

481       To illustrate the application of present PIML approach on screening of suitable

482       microbe-nutrient combination, a case study using synthetic data has been carried out.

483                                    [Table 3]

484     Table. 3 presents the microbial kinetic parameters for 4 different combinations of

485     microbes and nutrients, and rest all other parameters are kept constants. By feeding the

486     inputs through the trained neural network ML algorithm, the output oil recovery is

487     calculated and presented in the last column of Tab 3. It is inferred from Tab.3 that

488     among all the available combinations, the combination 4 shows highest oil recovery,

489     hence that corresponding microbe-nutrient pair can be used for field implementation.

490

491     *4. Conclusions*

492            Physics-Informed Machine Learning (PIML) approach is adopted to investigate

493     the performance of in-situ MEOR technique from limited experimental data, which is

494     difficult with conventional experimental and modelling approaches. Neural network ML

495     model used in the PIML approach had more accurately predicted the oil recovery. $Y_{PS}$,

496     flow velocity and initial oil saturation ($S_{ori}$) are correspondingly the most influential

497     microbial kinetic, operational and reservoir parameter. Higher oil recovery is achieved

498     by selecting a microbe-nutrient-reservoir pair having higher $Y_{PS}/Y_{XS}$ and $S_{ori}$ values.

499     This PIML approach helps to screen/identify suitable microbe-nutrient-reservoir pair at

500     initial laboratory stage itself, ensuring its success during the field implementation.

501

503     **References**

504  1. Ansah, E.O., Thanh, H.V., Sugai, Y., Nguele, R., Sasaki, K., 2020. Microbe-induced

505  fluid viscosity variation: field-scale simulation, sensitivity and geological uncertainty. J

506  Petrol. Explor. Prod. Technol. 10, 1983–2003.

507  2. Cruz, I.A., Chuenchart, W., Long, F., Surendra, K.C., Andrade, L.R.S., Bilal, M.,

508  Figueiredo, R.T., Khanal, S.K., Ferreira, L.F.R., 2021. Application of machine learning

509  in anaerobic digestion: Perspectives and challenges. Bioresour. Technol. 126433.

510  3. Joshi, S.J., Al-Wahaibi, Y.M., Al-Bahry, S.N., Elshafie, A.E., Al-Bemani, A.S., Al-

511  Bahry, A., Al-Mandhari, M.S., 2016. Production, characterization, and application of

512  bacillus licheniformis W16 biosurfactant in enhancing oil recovery. Front. Microbiol. 7,

513  1853.

514  4. Jeong, M.S., Lee, Y.W., Lee, H.S., Lee, K.S., 2021. Simulation-Based Optimization

515  of Microbial Enhanced Oil Recovery with a Model Integrating Temperature, Pressure,

516  and Salinity Effects. Energies, 14, 1131.

517  5. Jeong, M.S., Cho, J., Lee, K.S., 2022. Systematic modelling incorporating

518  temperature, pressure, and salinity effects on in-situ microbial selective plugging for

519  enhanced oil recovery in a multi-layered system, Biochem. Eng. J. 177, 108260.

520  6. Karniadakis, G.E., Kevrekidis, I.G., Lu, L., Perdikaris, P., Wang, S., Yang, L.,

521  2021. Physics-informed machine learning. Nat. Rev. Phys. 4, 422-440.

522  7. Keprate, A., Ratnayake, R.M.C., 2019. Data Mining for Estimating Fatigue Strength

523  Based on Composition and Process Parameters. Proc. ASME 2019 38th Int. Con. on

524  Ocean, Offshore and Arctic Eng. Vol 4: Materials Technology.

525  8. Liu, H., Zhang, J., Liang, F., Temizel, C., Basri, M.A., Mesdour, R., 2021.

526  Incorporation of physics into machine learning for production prediction from

527 unconventional reservoirs: a brief review of the gray-box approach. SPE Res. Eval.

528 Eng. 24, 847–858.

529 9. Markande, A.N., Patel, D., Varjani, S.J., 2021. A review on biosurfactants:

530 properties, applications and current developments. Bioresour. Technol. 330, 124963.

531 10. Nielsen, S.M., Nesterov, I., Shapiro, A.A., 2016. Microbial enhanced oil recovery-a

532 modeling study of the potential of spore-forming bacteria. Comput. Geosci. 20, 580.

533 11. Nikolova, C., Gutierrez, T., 2020. Use of microorganisms in the recovery of oil

534 from recalcitrant oil reservoirs: current state of knowledge, technological advances and

535 future prespective. Front. Microbiol. 10, 2996.

536 12. Sivasankar, P., Kanna, R., Kumar, G.S., Gummadi, S.N., 2016. Numerical

537 modelling of biophysicochemical effects on multispecies reactive transport in porous

538 media involving *Pseudomonas putida* for potential microbial enhanced oil recovery

539 application. Bioresour. Technol. 211, 348-359.

540 13. Sivasankar, P., Kumar, G.S., 2017. Influence of pH on dynamics of microbial

541 enhanced oil recovery processes using biosurfactant produced *Pseudomonas putida*:

542 Mathematical modelling and numerical simulation. Bioresour. Technol. 224, 498-508.

543 14. Shibulal, B., Al-Bahry, S.N., Al-Wahaibi, Y.M., Elshafie, A.E., Al-Bemani, A.S.,

544 Joshi, S.J., 2018. Microbial-Enhanced Heavy Oil Recovery under Laboratory

545 Conditions by *Bacillus firmus* BG4 and *Bacillus halodurans* BG5 Isolated from Heavy

546 Oil Fields. Colloids Interfaces. 2, 1.

547 15. Sivasankar, P., Kumar, G.S., 2019. Influence of bio-clogging induced formation

548 damage on performance of microbial enhanced oil recovery processes. Fuel. 236, 109.

549   16. Thanh, H.V., Sugai, Y., Sasaki, K., 2020. Application of artificial neural network

550   for predicting the performance of $CO_2$ enhanced oil recovery and storage in residual oil

551   zones. Sci. Rep. 10, 18204.

552   17. Tang, Q., Chen, Y., Yang, H., Liu, M., Xiao, H., Wang, S., Chen. H., Naqvi, S.R.,

553   2021. Machine learning prediction of pyrolytic gas yield and compositions with feature

554   reduction methods: Effects of pyrolysis conditions and biomass characteristics.

555   Bioresour. Technol. 339, 125581.

556   18. Varjani, S.J., Upasani, V.N., 2016. Core Flood study for enhanced oil recovery

557   through ex-situ bioaugmentation with thermo- and halo-tolerant rhamnolipid produced

558   by Pseudomonas aeruginosa NCIM 5514. Bioresour. Technol. 220, 175-182.

559   19. Varjani, S.J., Upasani, V.N., 2017. Critical review on biosurfactant analysis,

560   purification and characterization using rhamnolipid as a model biosurfactant. Bioresour.

561   Technol. 232, 389-397.

562   20. Wang, Z., Peng, X., Xia, A., Shah, A.A., Huang, Y., Zhu, X., Zhu, X., Liao, Q.,

563   2022. The role of machine learning to boost the bioenergy and biofuels conversion.

564   Bioresour. Technol. 343, 126099.

565   21. Zhang, W, Li, J., Liu, T., Leng, S., Yang, L., Peng, H., Jiang, S., Zhou, W., Leng,

566   L., Li, H., 2021. Machine learning prediction and optimization of bio-oil production

567   from hydrothermal liquefaction of algae. Bioresour. Technol. 342, 126011.

568

569

570

571

572

573

574

575

576

577

578

579

580

581

**Figures Caption**

1. Procedure of PIML approach followed in the present study to predict oil recovery by in-situ MEOR process

2. Frequency distribution of input and output data that are generated and used for training and testing of ML algorithms in the PIML approach

3. Validation of present microbial kinetic model results with measured experimental data for (a) variation of microbial concentration with time, (b) variation of sucrose concentration with time, (c) variation of biosurfactant concentration with time

590 4. (a) Pearson correlation and Spearman correlation coefficient matrix for microbe,

591 operational, reservoir and % oil recovery data, (b) Correlation coefficient values for

592 input microbe, operational and reservoir data towards output % oil recovery

593 5. (a) Relative Importance (RI) score of all input parameters, (b) RI score of input

594 microbial-nutrient parameters, (c) RI score of input operational parameters, (d) RI score

595 of input reservoir parameter in predicting output % of oil recovery

596 6. Comparative performance of 12 different ML algorithms in predicting the oil

597 recovery against the actual % of oil recovery for in-situ MEOR application.

598 **Tables Caption**

599 1. Input parameters and their corresponding value range used in the present study

600 2. Performance of different ML algorithms in predicting the actual % of oil recovery

601 3. Microbial kinetic parameters for different microbial-nutrient combinations and their

602 corresponding oil recovery determined using PIML modelling approach.

603 Table 1. Input parameters and their corresponding value range used in the present study

604 [Sivasankar et al. (2016)]]

| Parameter | Reference value | Range |
|---|---|---|
| $Y_{XS}$ | 0.1843 | 0.092 - 0.276 |
| $Y_{PS}$ | 0.078 [Sivasankar et al. (2016)] | 0.03900733 - 0.116996715 |
| $K_{XS}$ (*g/l*) | 6.86 [Sivasankar et al. (2016)] | 3.430058808 - 10.28959695 |
| $U_{max}$ (*h^-1*) | 0.053 [Sivasankar et al. (2016)] | 0.02650387 - 0.079495509 |

| $X_i$ (g/l) | 0.1521167 [Sivasankar et al. (2016)] | 0.076094593 - 0.22824074 |
|---|---|---|
| $S_i$ (g/l) | 19.234 [Sivasankar et al. (2016)] | 9.617601084 - 28.84936382 |
| $A_i$ (g/l) | 3 [Sivasankar et al. (2016)] | 1.500165456 - 4.49988547 |
| $T_r$ (h) | 150 | 100 – 200 |
| $u_w$ (m/h) | 0.0004 [Sivasankar et al. (2016)] | 0.0002 – 0.0006 |
| $\mu_w$ (Nhm$^{-2}$) | 0.001 [Sivasankar et al. (2016)] | 0.0005 – 0.0015 |
| Initial IFT (mN/m) | 51.6 [Sivasankar et al. (2016)] | 25.80405697 - 77.39695 |
| $S_{wir}$ | 0.2 | 0.10000517 - 0.299989777 |
| $S_{ori}$ | 0.4 | 0.20000454 - 0.599986979 |
| Output - % oil recovery | Mean - 7.423469637 Median - 5.510703244 | 0.174742622 - 48.23409386 |

605
606
607

608   Table 2: Performance of different ML algorithms in predicting the actual % of oil

609   recovery

| **Model** | $R^2$ | *RMSE* | *Explained Variance Score* |
|---|---|---|---|
| KNN | 0.369897442 | 3.369742385 | 0.370031425 |
| Decision Trees | 0.408587985 | 4.356687642 | 0.408683171 |
| Lasso | 0.510880307 | 3.697724477 | 0.51138261 |
| Ridge | 0.512474691 | 3.697351147 | 0.512976983 |

| | | | |
|---|---|---|---|
| Linear Regression | 0.51299017 | 3.697264614 | 0.513491055 |
| Random Forests | 0.639724278 | 2.941482946 | 0.63972701 |
| ADA Boost | 0.639746354 | 2.811959461 | 0.639980378 |
| Gradient Boosting | 0.896025654 | 1.851540808 | 0.896025658 |
| Gaussian Process | 0.951787039 | 1.353896887 | 0.951802595 |
| Polynomial (4) | 0.963022723 | 1.26029449 | 0.963039744 |
| SVR | 0.964436929 | 1.184682456 | 0.964596291 |
| Neural Network | 0.987349995 | 0.714597767 | 0.987656577 |

Table 3: Microbial kinetic parameters for different microbial-nutrient combinations and

their corresponding oil recovery determined using PIML modelling approach

| Combinations | $Y_{XS}$ | $Y_{PS}$ | $K_{XS}$ | $U_{max}$ | Output Oil Recovery, % |
|---|---|---|---|---|---|
| 1 | 0.098734 | 0.067978 | 4.077158 | 0.068247 | 6.3529167 |
| 2 | 0.148067 | 0.081408 | 4.763728 | 0.056552 | 4.531987 |

| | | | | | |
|---|---|---|---|---|---|
| 3 | 0.169445 | 0.091877 | 4.923773 | 0.053963 | 4.549041 |
| 4 | 0.121985 | 0.095722 | 8.262717 | 0.038336 | 7.3623314 |

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

| Classification and collection of input parametric data (*Refer Sec. 2.1*) |

| Generation of large input data sets (*Refer Sec. 2.2.1*) <br> For all identified input parameters, a certain % of SD is considered to their corresponding collected value and obtained a value range. Within this value range, large data set (10000 in present study) of input values are generated. |

| Generation of large output data sets from physics based model (*Refer Sec. 2.2.2*) <br> Using the generated input data, output oil recovery data is generated by simulation of physics-based model. Large data sets (10000 in the present study) of output data are generated by performing 10000 simulations of physics-based model. |

| Quantification of interaction strength between input and output parameters (*Refer Sec. 2.3*) <br> By statistically analysing the generated input and output data sets, PCC and SPCC values are determined which quantifies the interaction strength between input and output parameters. |

| Splitting of generated input and output data in 7:3 ratio for training & testing of 12 different ML algorithms in PIML approach |

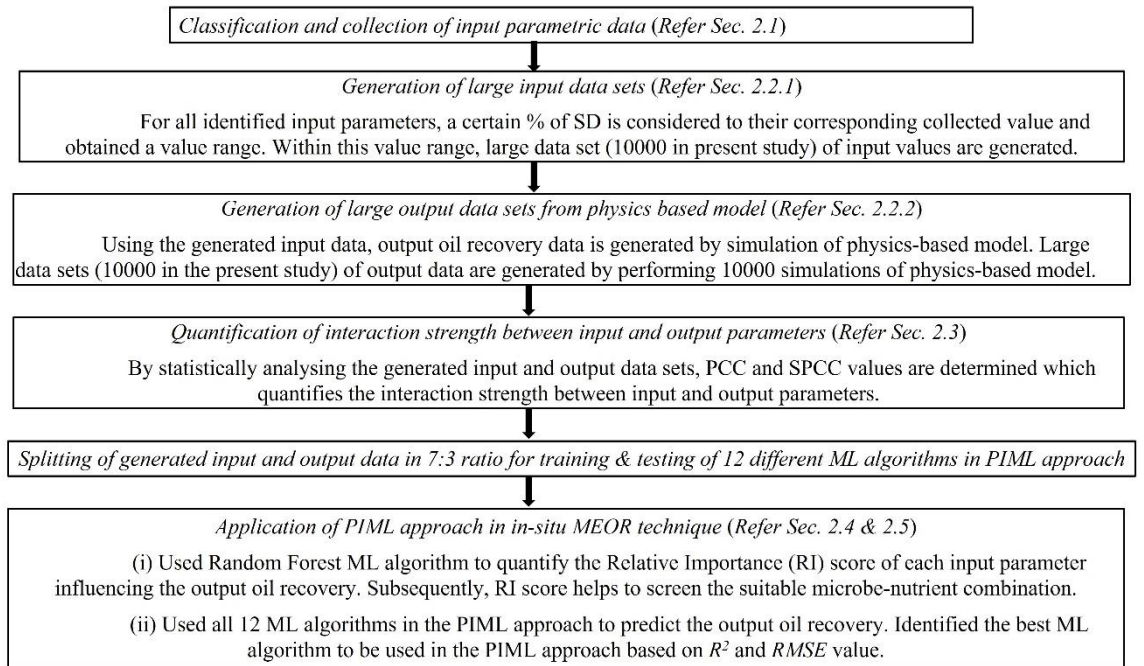| Application of PIML approach in in-situ MEOR technique (*Refer Sec. 2.4 & 2.5*) <br> (i) Used Random Forest ML algorithm to quantify the Relative Importance (RI) score of each input parameter influencing the output oil recovery. Subsequently, RI score helps to screen the suitable microbe-nutrient combination. <br> (ii) Used all 12 ML algorithms in the PIML approach to predict the output oil recovery. Identified the best ML algorithm to be used in the PIML approach based on $R^2$ and *RMSE* value. |

637

638 **Figure 1**: Procedure of PIML approach followed in the present study to predict oil

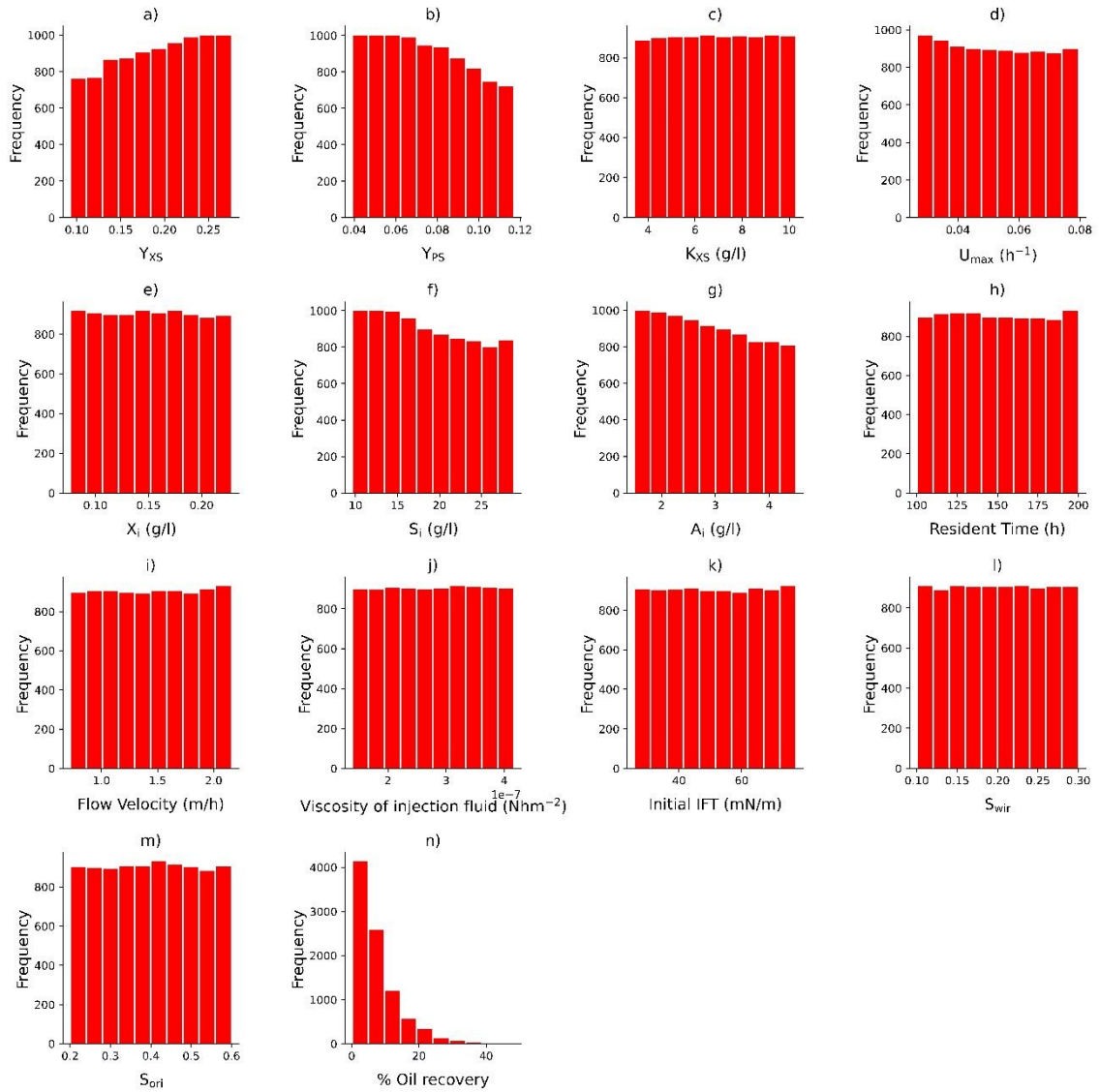639 recovery by in-situ MEOR process (NIKHIL WILL CHANGE IT)

640

641

642

643

644

645

646

647

648

649

650

**Figure 2**: HIstogram of input and output data that are generated and used for training and testing of ML algorithms in the PIML approach.

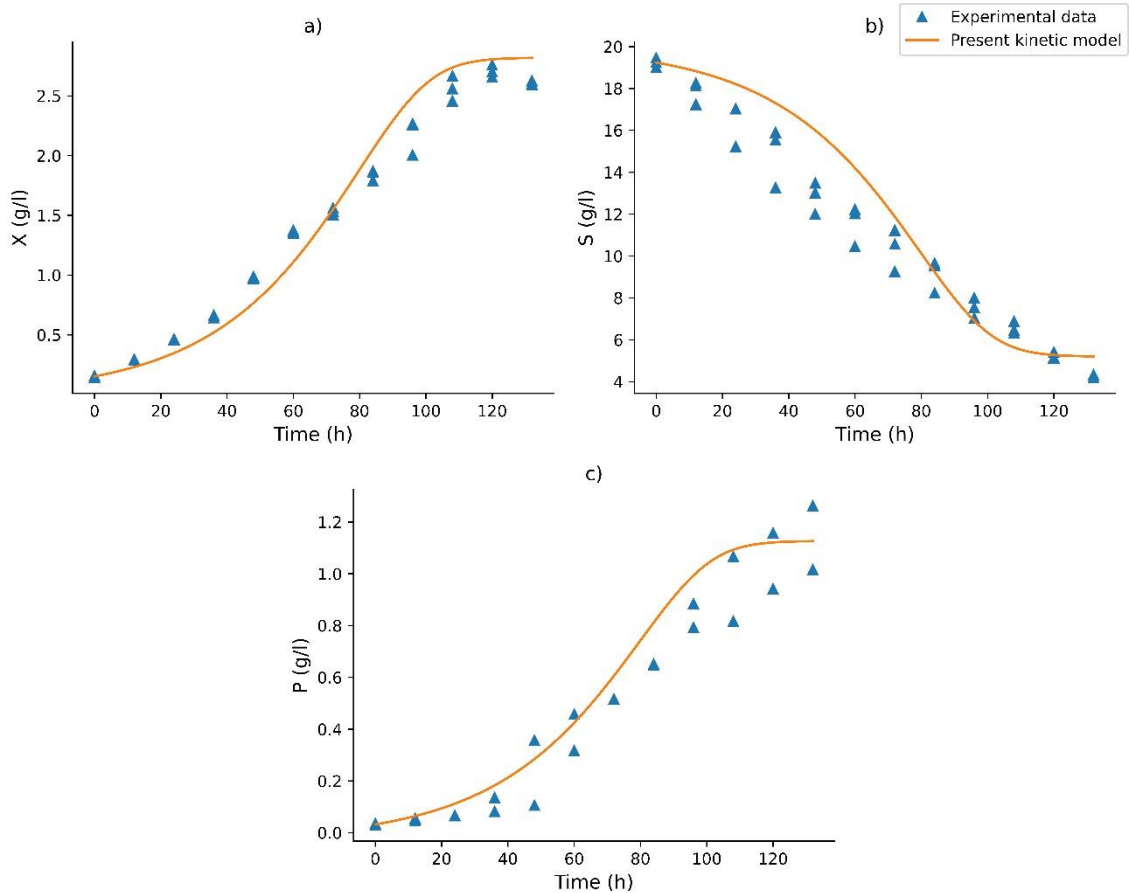651

652

653

654

655

656

657

**Figure 3**: Validation of present microbial kinetic model results with measured experimental data for (a) variation of microbial concentration with time, (b) variation of sucrose concentration with time, (c) variation of biosurfactant concentration with time.
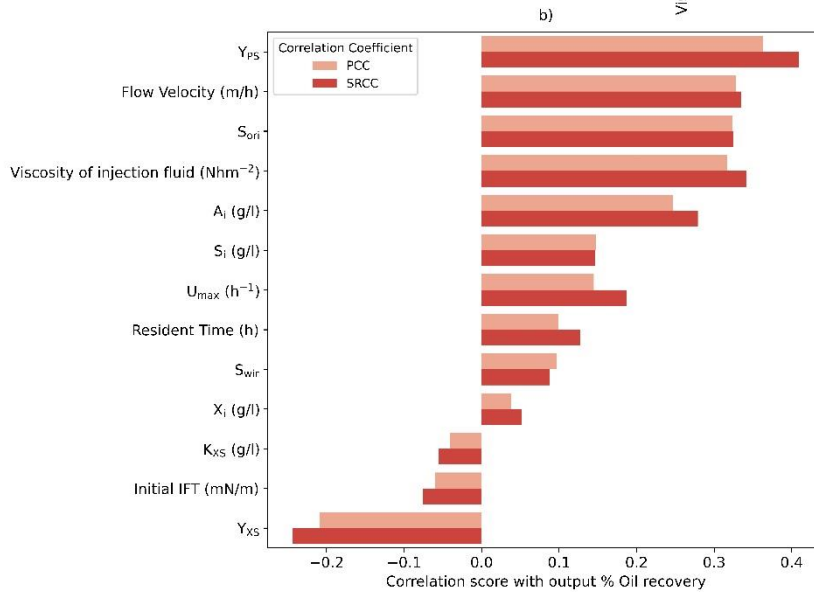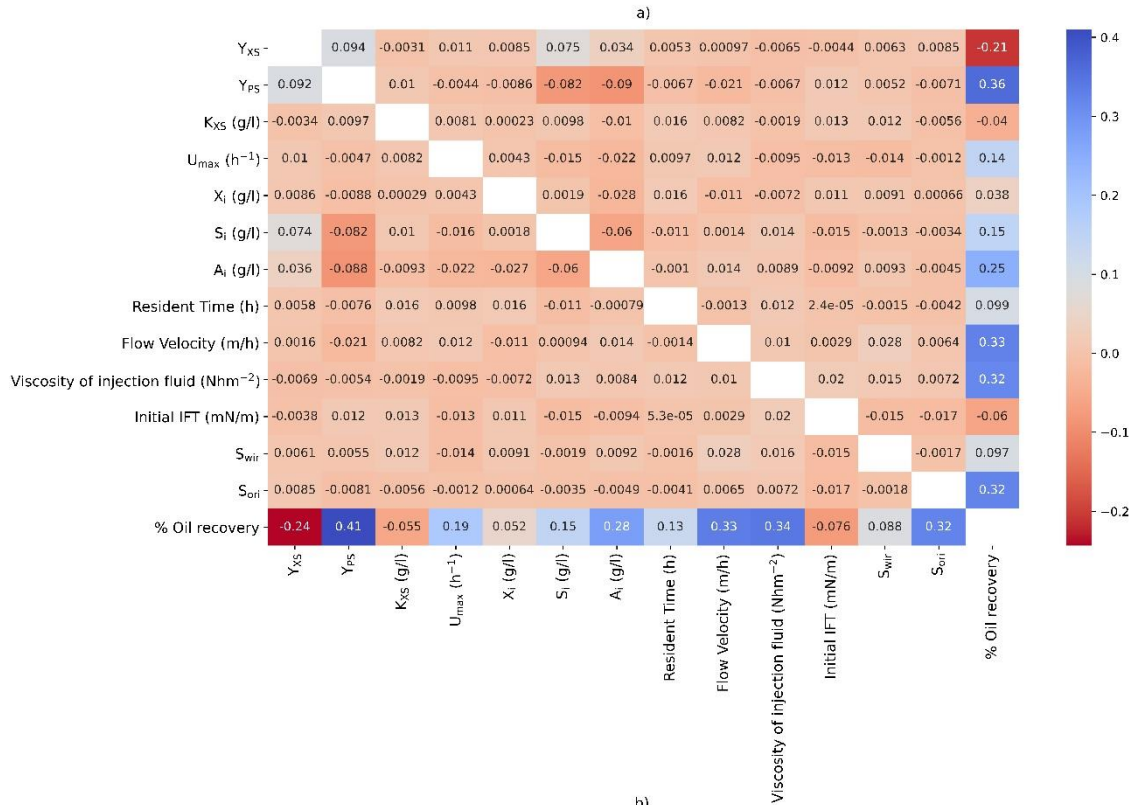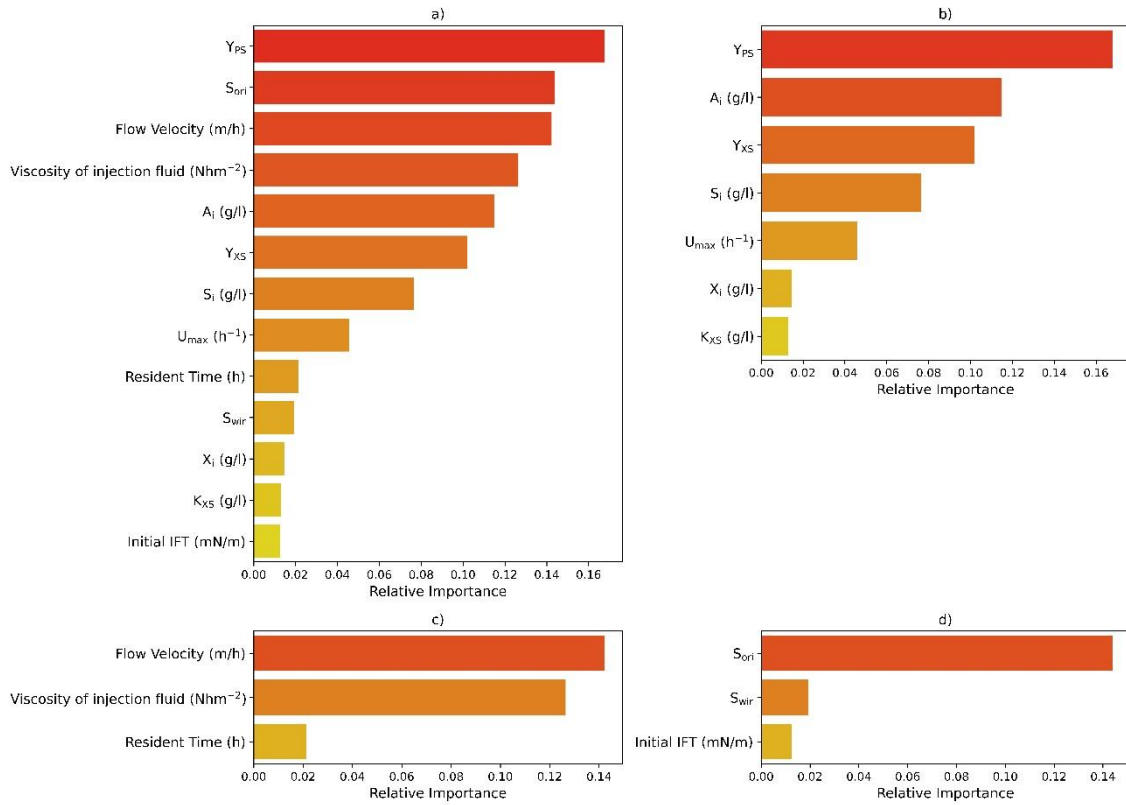
**Figure 4:** (a) Pearson correlation and Spearman correlation coefficient matrix for microbe, operational, reservoir and % oil recovery data, (b) Correlation coefficient values for input microbe, operational and reservoir data towards output % oil recovery.

**Figure 5:** (a) Relative Importance (RI) score of all input parameters, (b) RI score of input microbial-nutrient parameters, (c) RI score of input operational parameters, (d) RI score of input reservoir parameter in predicting output % of oil recovery.
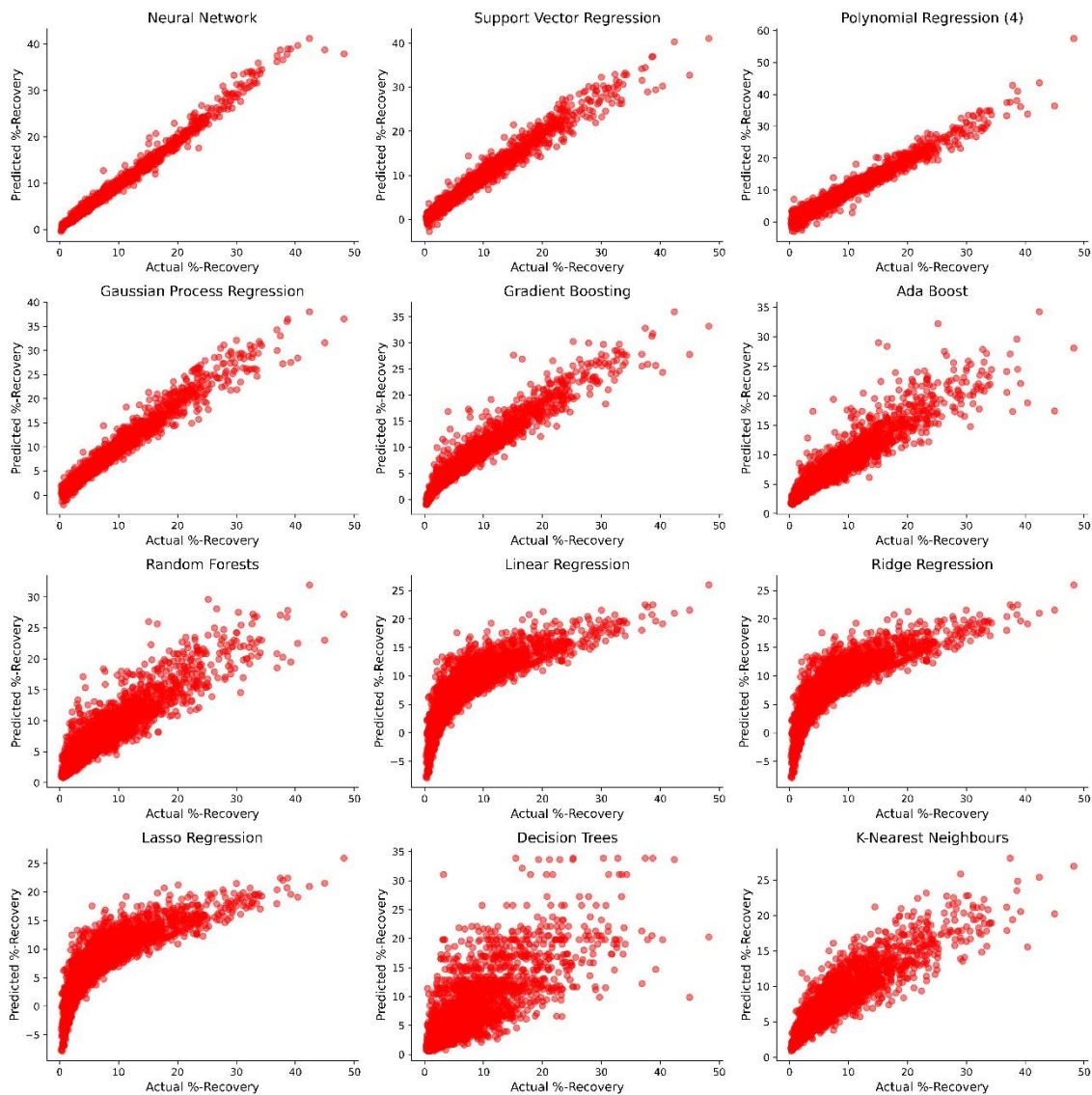
684

**Figure 6:** Comparative performance of 12 different ML algorithms in predicting the oil

recovery against the actual % of oil recovery for in-situ MEOR application.

687