

Energy Advances

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: K. Sandunil, Z. Bennour, H. Ben Mahmud and A. Giwelli, *Energy Adv.*, 2024, DOI: 10.1039/D4YA00313F.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

Effects of Tuning Decision Trees in Random Forest Regression on Predicting Porosity of a Hydrocarbon Reservoir. Case Study: Volve Oil Field, North Sea

Kushan Sandunil^{a*}, Ziad Bennour^b, Hisham Ben Mahmud^c, Ausama Giwelli^{d,e}

^{a,b} Curtin University Malaysia, 98009 Miri, Sarawak, Malaysia

^c Universiti Teknologi PETRONAS, 32610 Seri Iskandar, Perak Darul Ridzuan, Malaysia

^d INPEX, 100 St Georges Terrace 6000 Perth WA, Australia

^e WASM, Curtin University, Kensington WA 6151, Australia

*Corresponding author: kwkushan@postgrad.curtin.edu.my

Abstract: Machine learning (ML) has emerged as a powerful tool in petroleum engineering for automatically interpreting well logs and characterizing reservoir properties such as porosity. As a result, researchers are trying to enhance the performance of ML models further to widen their applicability in the real world. Random forest regression (RFR) is one such widely used ML technique which was developed by combining multiple decision trees. To improve its performance one of its hyperparameters, the number of trees in the forest ($n_{estimators}$) are tuned during model optimization. However, existing literature lacks in-depth studies on the influence of $n_{estimators}$ on RFR model when used for predicting porosity, given that $n_{estimators}$ is one of the most influential hyperparameters that can be tuned to optimize the RFR algorithm. In this study, the effects of $n_{estimators}$ on RFR model in porosity prediction were investigated. Furthermore, $n_{estimators}$ ' interaction with another two key hyperparameters, namely the number of features considered for the best split ($max_features$) and the minimum number of samples required to be at a leaf node ($min_samples_leaf$) was explored. The RFR models were developed using 4 input



features namely, resistivity log (RES), neutron porosity log (NPHI), gamma ray log (GR) and corresponding depths obtained from Volve oil field in North Sea and calculated porosity was used as the target data. The methodology consisted of 4 approaches. In the first approach only $n_estimators$ were changed, in the second approach $n_estimators$ were changed along with $max_features$, in the third approach $n_estimators$ were changed along with $min_samples_leaf$ and in the final approach all three hyperparameters were tuned. Altogether 24 RFR models were developed, and models were evaluated using adjusted R^2 (adj. R^2), root mean squared error (RMSE) and their computational times. The obtained results showed that the highest performance with an adj. R^2 value of 0.8505 was given when $n_estimators$ was 81, $max_features$ was 2 and $min_samples_leaf$ was 1. In approach 2, when $n_estimators$ upper limit was increased from 10 to 100 there was a test model performance growth of more than 1.60%, whereas increasing $n_estimators$ ' upper limit from 100 to 1000 showed a performance drop of around 0.4%. Models developed by tuning $n_estimators$ from 1 to 100 in intervals of 10 had healthy test model adj. R^2 values and lower computational times making them the best $n_estimators$ range and interval when both performances and computational times were taken into consideration to predict porosity of Volve oil field in North Sea. Further, it was concluded that by tuning only $n_estimators$ and $max_features$ the performance of RFR models can be increased significantly.

Keywords: *Machine Learning, Random Forest Regression, Porosity Prediction, Hyperparameter Tuning, Decision Trees Tuning, Volve Oil Field*

1. Introduction

Artificial intelligence (AI) has become a popular topic over the past few years due to its immense potential in STEM fields. Machine learning (ML) is a branch of AI where it learns with or without supervision to do predictions. Its ability to predict or forecast outputs, decrease computational time



and extract features from complex and high-dimensional datasets make it a tremendous tool to work with complex and huge datasets.¹⁻³ The concept of ML was first put forward by Turing.⁴ Since then, ML has seen a significant improvement with the invention of complex and high performing algorithms. With the popularity of ML algorithms, like many other engineering sectors, their applicability in reservoir engineering has been tested, especially in porosity prediction. Porosity gives an idea about the fluid storage capacity, and it plays a vital role in the upstream oil and gas industry, since it is used in estimating petroleum initially in place in the reservoir. Core analysis is a reliable and a widely accepted approach used to estimate porosity. However, this method is expensive and time consuming. To address these challenges petroleum engineers and researchers are investigating the applicability of ML in reservoir characterization. Random forest regression (RFR) is one such ML algorithm which has successfully been used to predict porosity. To enhance the performance of ML models hyperparameter optimization is used. The research investigated the effects of one of the main hyperparameters of RFR, number of decision trees in the forest ($n_estimators$) when predicting porosity of a sandstone dominated section in Volve oilfield. Moreover, the behaviour of two other widely used hyperparameters in RFR, minimum number of samples required to be at a leaf node ($min_samples_leaf$) and number of features considered for the best split ($max_features$) when tuned along with $n_estimators$ was studied. Apart from the primary objectives mentioned, this study tested the feasibility of an optimized RFR algorithm in reservoir characterization, specifically porosity prediction which could be further extended to be used in permeability and saturation prediction in future studies.

ML application in reservoir characterization has seen a significant increase over the last couple of decades due to its ability to tackle regression and classification type problems.⁵⁻⁷ With the evolution of ML, a notable number of algorithms have been introduced. Artificial neural network



(ANN) which uses a parallel processing approach and developed based on the function of a neuron of a human brain has been utilized in petrophysical parameter prediction.^{8,9} Support vector regression (SVR) is another algorithm developed in the initial stages of the ML timeline and has the capability to handle non-linear relationships between a set of inputs and an output. Moreover, SVR has been utilized widely in reservoir characterization.¹⁰⁻¹³ Least absolute shrinkage and selection operator (LASSO) regression and Bayesian model averaging (BMA) has also been used in ML related studies extensively in the literature.¹⁴ BMA uses Bayes theorem and LASSO uses residual sums of squares to build a linear relationship between the inputs and the output. BMA and LASSO regression has been used in permeability modelling in recent studies.⁵ Apart from petrophysical parameter prediction ML models have been used in lithofacies classification as well.¹⁵ Generally, these studies utilized ML approaches to model lithofacies sequences as a function of well logging data in order to predict discrete lithofacies distribution at missing intervals.¹⁶⁻¹⁸ Besides permeability prediction, water saturation estimation and lithofacies classification, ML models have been used in reservoir porosity estimation which is the parameter focused on this study. ML algorithms such as ANN, deep learning and SVR were used to predict porosity using logging data, seismic attributes and drilling parameters.¹⁹⁻²¹

Apart from the aforementioned ML models an ML approach known as ensemble learning has been applied in many recent studies. Here, ML base models (weaker models) are strategically combined to produce a high performing and efficient model as shown in Fig. 1. Ensemble ML models has become a popular tool among researchers to predict petrophysical properties due to their ability to reduce overfitting and underfitting.²²⁻²⁶ RFR is one such popular ensemble ML model which was developed by amalgamating multiple decision trees.²⁷



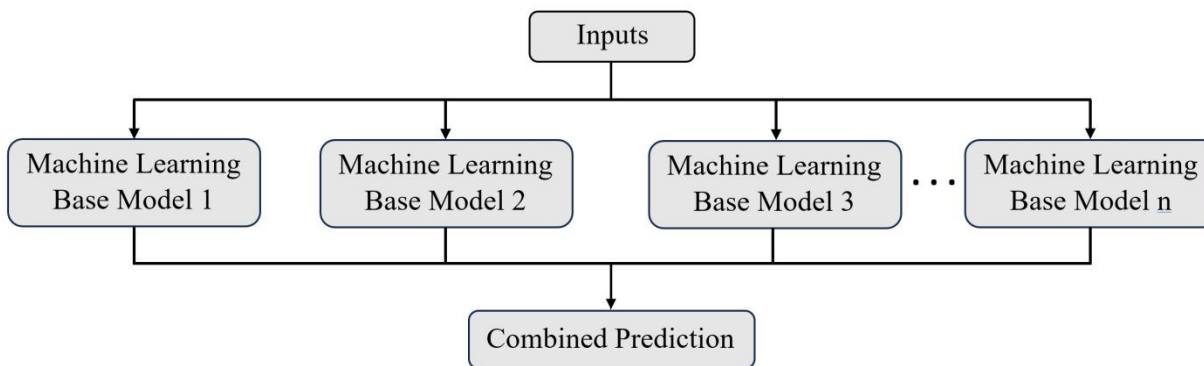


Fig. 1 Representation of the ensemble model.

Hyperparameter tuning is a process implemented to fine-tune ML algorithms to obtain optimal models.²⁸⁻³⁰ There are several hyperparameters that can be controlled in an RFR model such as, $n_estimators$, $max_features$, $min_samples_leaf$, maximum depth of the tree (max_depth), fraction of the original dataset assigned to any individual tree ($max_samples$), minimum number of samples required to split an internal node ($min_samples_split$), maximum leaf nodes to restrict the growth of the tree (max_leaf_nodes).

Hyperparameter optimization has been utilized in recent studies related to reservoir characterization. Wang et al. developed a RFR model to predict permeability in Xishan Coalfield, China.²⁴ Five hyperparameters $n_estimators$, $max_features$, max_depth , $min_samples_leaf$ and $min_samples_split$ were tuned during hyperparameter optimization. Zou et al. estimated reservoir porosity using a random forest algorithm.³¹ During the hyperparameter optimization stage $n_estimators$, $max_features$, $min_samples_leaf$, $min_samples_split$ and max_depth were tuned. Rezaee and Ekundayo tuned $n_estimator$, $min_samples_leaf$, $min_samples_split$ and max_depth during the development of the RFR model used to predict permeability of precipice sandstone in Surat Basin, Australia.³²



Even though hyperparameters have been tuned during hyperparameter optimization phase of an ensemble ML model development, literature lacks studies done specifically focusing on the effects of hyperparameter tuning in ensemble learning when predicting petrophysical properties in reservoir characterization. Addressing this research gap, in this study, the authors investigated the influence of one of the most utilized hyperparameters in the literatures, *n_estimators* of RFR when predicting porosity of a hydrocarbon reservoir. Also, the effects of *n_estimators* were studied along with another two widely used hyperparameters, *max_features* and *min_samples_leaf* when predicting porosity of Volve oil field in North Sea. The study considered a supervised learning regression approach. The workflow of the study consisted of data preprocessing, RFR models development and models analysis. Several RFR models were developed tuning *n_estimators*, tuning *n_estimators* along with *max_features*, tuning *n_estimators* along with *min_samples_leaf* and tuning all three hyperparameters at once under four approaches by integrating grid search optimization and K-fold cross-validation. Models' performances were evaluated based on adjusted coefficient of determination (adj. R^2), root mean squared error (RMSE) and computational time. The study considered only aforementioned 3 hyperparameters due to processing capacity limitations. However, the study is expected to be a solid initiation towards the development of future studies on the effects of hyperparameters in ML algorithms in reservoir characterization.

2. Methodology

Developing an ML model consists of multiple steps, namely, data acquisition, data preprocessing models development and data analysis.³³⁻³⁵ In this study the above steps were implemented to develop robust ML models.

2.1 Geological setting and dataset



Volve oil field (Fig. 2) was selected as the study area and the well log data of the field was publicly available. Several ML related studies have been conducted using the Volve oil field datasets.³⁶⁻³⁸ It was formed during the Jurassic period by the collapse of adjacent salt ridges. Oil was discovered in the field back in 1993 in the middle Jurassic Hugin sandstone formations identifying it as a clastic reservoir.

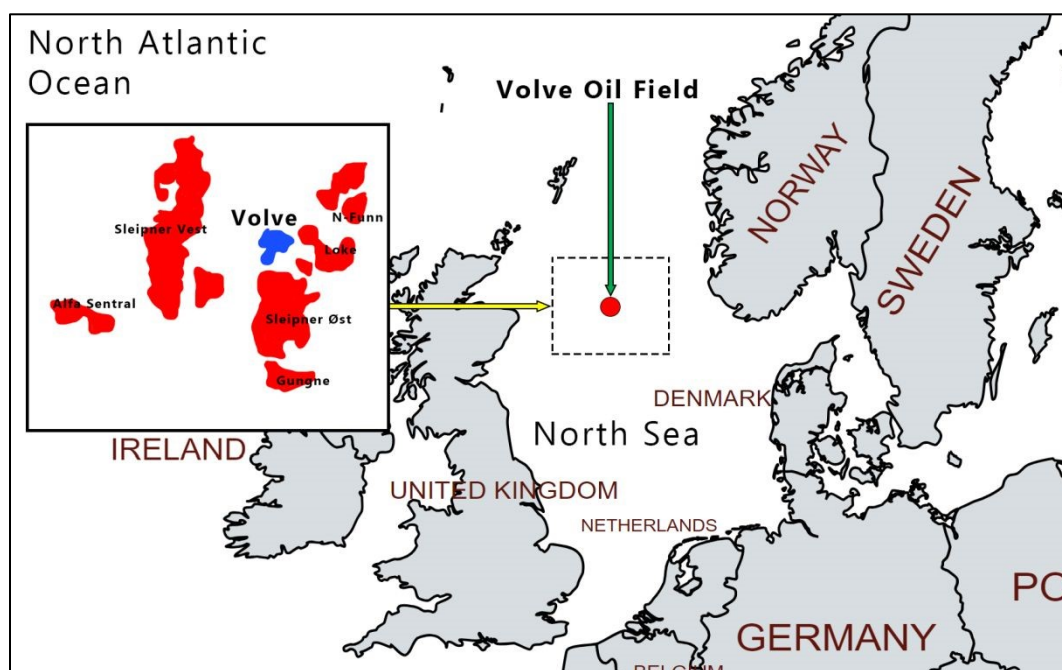


Fig. 2 Study area – Volve oil field’s location in the North Sea. Adapted from Mapchart.³⁹

The Hugin Formation is 153 m thick and oil-bearing, and was penetrated at 3796.5 m, approximately 60 m deeper than expected. The total oil column in the well was 80 m, but no clear oil-water contact was observed.^{38,40} The reservoir section was made up of highly variable fine to coarse grained, well to poorly sorted subarkosic arenite sandstones with good to excellent reservoir properties. The Hugin Formation of the area consists of shallow marine shoreface, coastal plain/lagoonal, channel and possibly mouth bar deposits. The underlying Skagerrak formation was completely tight due to extensive kaolinite and dolomite cementation. The current study used data



from the well 15/9-19A. The well was drilled through the Skagerrak formation and terminated approximately 30 m into the Triassic Smith Bank formation. To fully utilize the available data, the study considered data from 3666.59 to 3907.08 m depth interval. This depth interval ran through three formations namely, Draupne, Heather and Hugin. The stratigraphic column and description about the vertical facies distribution of the focused section is shown in Fig. 3.



| Stratigraphy | | | | Depth (m) | Lithology | Lithological Description | |
|--------------|----------------|----------------|-------------------|-----------|--|---|--|
| System | Series | Group | Formation | | | | |
| Jurassic | Upper Jurassic | Viking Group | Draupne Formation | 3666.5 | 3700 | <p>Claystone: Dark brown - brownish black, firm, blocky - subfissile, uniform, non-calcareous, carbonaceous</p> <p>Claystone: Predominantly dark brown - brownish black - dusky yellowish brown, firm, blocky - subfissile, silty, non-calcareous, micromicaceous, carbonaceous</p> | |
| | | | Heather Formation | 3706.0 | | 3750 | <p>Claystone: Olive black - grey black, firm, blocky, silty and sandy in parts</p> <p>Limestone: Medium grey - off white firm, blocky, silty and sandy in parts</p> <p>Claystone: Olive black, firm - moderately hard, blocky silty, carbonaceous, micromicaceous, non-calcareous, rarely medium grey - dark green grey</p> <p>Claystone: Medium dark grey - dark grey, occasionally green grey, firm</p> <p>Limestone: Medium dark grey - off white yellow brown, firm - moderate hard</p> |
| | | Vestland Group | Hugin Formation | | 3796.5 | | 3800 |
| | | | | | 3850 | 3900 | |
| | | | 3908.0 | | <p>Dolomite: Light grey brown, moderate hard, angular trace macrofossils, tight</p> | | |

Fig. 3 Stratigraphic column and facies description of the considered subsurface section. Adapted from Statoil.⁴¹

The dataset consisted of depth, well log data and corresponding calculated porosity values and had a total of 1547 data points. Three well log parameters, namely: resistivity log (RES), neutron



porosity log (NPHI) and gamma ray log (GR) along with corresponding depth were used as input features and total porosity (PHIF) was used as the target data. PHIF was calculated using porosity from density log (PHID) and NPHI. PHIF was derived from the density log which is calibrated to overburden corrected core porosity for wells drilled with either oil-based mud or water-based mud. NPHI was used to correct for varying mud filtrate invasion. Equations used to calculate PHIF and PHID are shown in Eq. 1 and Eq. 2 respectively.

$$PHIF = PHID + A \times (NPHI - PHID) + B. \quad (1)$$

$$PHID = \frac{\rho_{ma} - \rho_b}{\rho_{ma} - \rho_{fl}} \quad (2)$$

In Eq. 1 A and B are regression coefficients and in Eq. 2 ρ_{ma} is matrix density, ρ_b is measured bulk density and ρ_{fl} is pore fluid density. Calculated PHIF values were assumed as actual porosities during model development and evaluation.

2.2 Data preprocessing

The raw data acquired from Volve oil field was subjected to data preprocessing before they were used in ML model development. Three main data preprocessing practices; (i) data cleaning, (ii) feature scaling and (iii) data division were utilized in this study.⁴²⁻⁴⁹ Under data cleaning, missing values and outliers were identified. Missing values were the sections where datapoints were missing from the dataset. Outliers were the data which lied outside of a considered range in each feature. The interquartile range method was used to detect the outliers. Both missing data and outliers were treated by removing them completely from the dataset.^{50,51}

Feature scaling is also a common practice implemented during data preprocessing. There are two widely used feature scaling approaches in the literature namely, normalization and standardization.



However, in this study feature scaling was neglected since RFR is a tree-based ML model where splits do not change with any monotonic transformation.⁵²

Data division was carried out by splitting the dataset into 2 parts as training and testing. Training portion was used to train the ML models while the testing portion was used to test the trained models. Train-test ratio was considered as 80:20, i.e., 80% of the total dataset was allocated for training while the remaining 20% was used for testing.^{53,54}

2.3 Machine learning models development

The RFR model is a combination of multiple decision trees. A typical architecture of an RFR model is shown in Fig. 4. Segal demonstrated random forest algorithm mathematically as $h(x; \theta_r)$, $r = 1, \dots, R$ where x represents the observed input vector associated with vector X . X and θ_r are independent and identically distributed random vectors.⁵⁵ For mathematical clarification let's define a vector with numerical outcomes Y . Therefore, training dataset of the RFR can be assumed to be drawn from a joint distribution of (X, Y) .

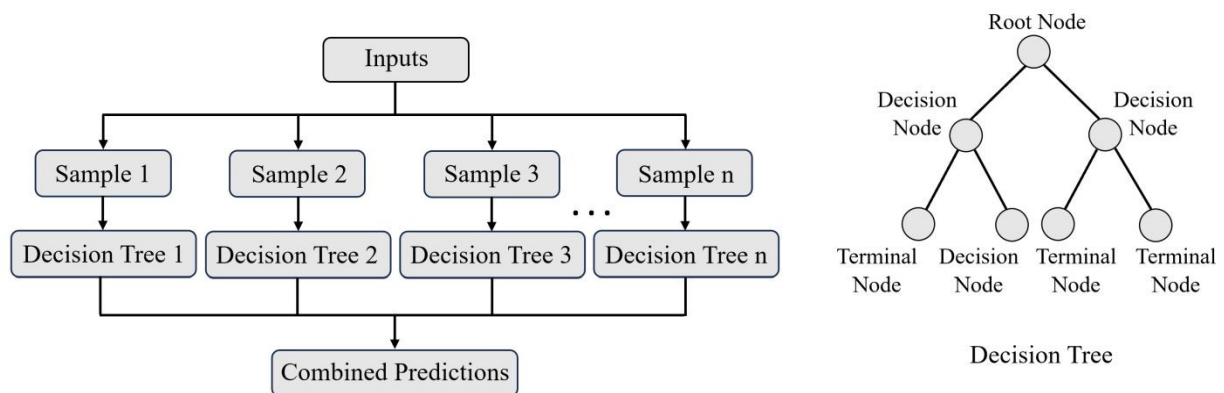


Fig. 4 Random Forest architecture (left) and the base model architecture (right).

For regression, the random forest prediction is the unweighted average over the collection:



$$h(x) = \left(\frac{1}{R}\right) \sum_{r=1}^R h(x; \theta_r). \quad (3)$$

As $r \rightarrow \infty$ the Law of Large Numbers ensures,

$$E_{X,Y}(Y - \bar{h}(X))^2 \rightarrow E_{X,Y}(Y - E_{\theta}h(X; \theta))^2. \quad (4)$$

The quantity on the right is the prediction (or generalization) error for the random forest, designated PE_f^* . The convergence in Eq. (4) implies that random forests do not overfit.

Now define the average prediction error for an individual tree as $h(X; \theta)$

$$PE_f^* = E_{\theta} E_{X,Y}(Y - h(X; \theta))^2. \quad (5)$$

Assuming that for all θ the tree is unbiased, i.e., $EY = E_X h(X; \theta)$. Then

$$PE_f^* \leq \bar{\rho} PE_f^* \quad (6)$$

Where $\bar{\rho}$ is the weighted correlation between residuals $Y - h(X; \theta)$ and $Y - h(X; \theta')$ for independent θ, θ' .

The inequality shown by Eq. (6) highlights what is required for accurate RFR which is having a low correlation between residuals of differing tree members of the forest and low prediction error for the individual trees. Model's performance can be further enhanced by tuning its hyperparameters.

During the study RFR models were developed using Python programming language. The cleaned dataset obtained during data preprocessing stage was loaded into Python. Then they were split into training and testing. Python based *scikit-learn* library's *RandomForestRegressor* was used to develop the RFR algorithm. The *RandomForestRegressor* comes with default hyperparameters



built into it. Default values assigned to some of the main hyperparameters of RFR in *scikit-learn* are tabulated in Table 1.

Table 1 Some hyperparameters of random forest algorithm and their default values in *scikit-learn* library.

| Hyperparameter | Default Value |
|--------------------------|---------------|
| <i>n_estimators</i> | 100 |
| <i>max_features</i> | 1.0 |
| <i>min_samples_leaf</i> | 1 |
| <i>max_depth</i> | None |
| <i>max_samples</i> | None |
| <i>min_samples_split</i> | 2 |
| <i>max_leaf_nodes</i> | None |

However, rather than using the default hyperparameters assigned by *scikit-learn* library, to achieve the primary objectives of the study hyperparameter optimization was implemented. Hyperparameter optimization is a commonly used practice to build robust ML models.^{56,57} The hyperparameters of RFR were tuned using the grid search optimization (GSO) approach. For this *GridSearchCV* optimization algorithm in *scikit-learn* library was used. GSO was considered since it runs through all the possible combinations in the hyperparameter space hence selecting the best combination of the space.^{57,58} The hyperparameter space was predefined by including the possible values and it was fed into the GSO algorithm.



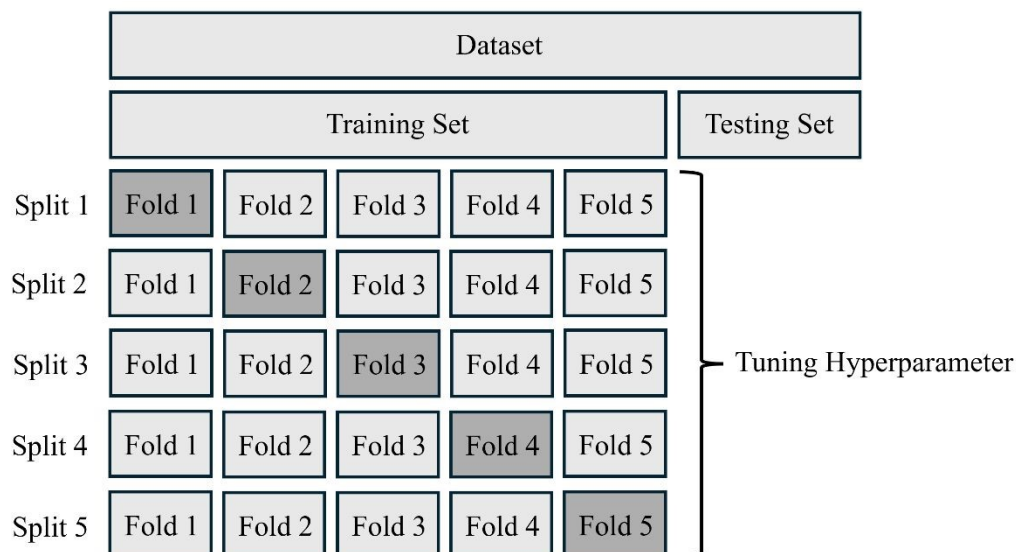


Fig. 5 Demonstration of K-fold cross validation.

GSO was implemented along with random subsampling cross-validation. An approach known as K-fold cross-validation was used. During K-fold cross-validation the training dataset is divided into K number of same-sized portions (folds) and K-1 of the portions will be used for training and the remainder will be used for validation.^{59,60} This will be repeated until each fold gets the chance to be the validation set. For this study 5-fold cross-validation was implemented as shown in Fig. 5. Therefore, the training set was divided into five portions and during each split four folds were used for training and one fold was used for validation.



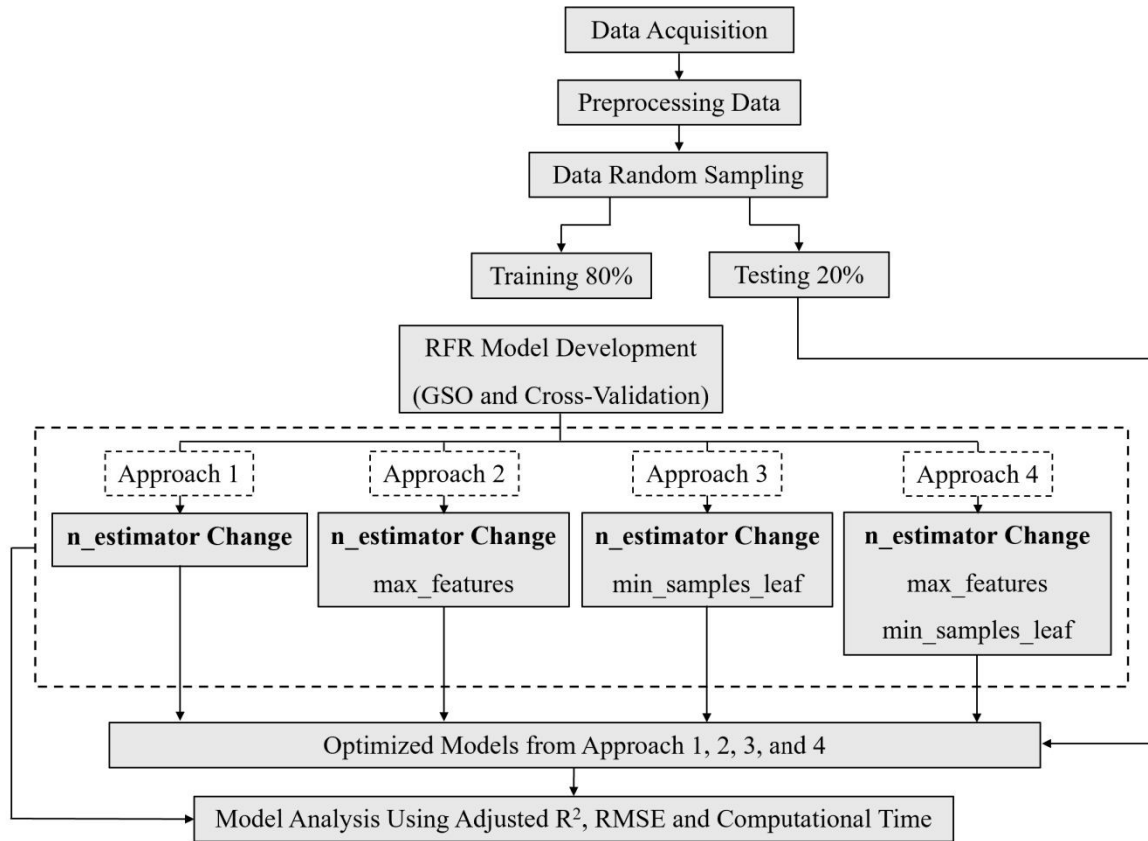


Fig. 6 Workflow of the methodology.

Tuning was done under 4 approaches as shown in Fig. 6 to investigate the effects of the considered hyperparameters. In the first approach $n_estimators$ was changed from 1 to 10, 1 to 100 and 1 to 1000 in different intervals. The notation used to demonstrate the $n_estimators$ change is shown in Table 2.



Table 2 Notations of $n_estimators$ change and their representations.

| $n_estimators$ Change | Starting Value | Ending Value | Increment |
|------------------------|----------------|--------------|-----------|
| Notation | | | |
| 1:10:1 | 1 | 10 | 1 |
| 1:100:1 | 1 | 100 | 1 |
| 1:100:10 | 1 | 100 | 10 |
| 1:1000:1 | 1 | 1000 | 1 |
| 1:1000:10 | 1 | 1000 | 10 |
| 1:1000:100 | 1 | 1000 | 100 |

In the second approach $n_estimators$ was changed from 1 to 1000 in the same way as approach 1 along with $max_features$. Here, $max_features$ was changed from 10% to 100% of total features in increments of 10%. In the third approach $n_estimators$ was changed in the same way along with $min_samples_leaf$. In this case, $min_samples_leaf$ was changed from 1 to 20 in intervals of 1. In the fourth approach all 3 hyperparameters, i.e., $n_estimators$, $max_features$ and $min_samples_leaf$ were varied at the same time in above mentioned intervals. In each approach values of all the other hyperparameters of RFR were kept at their default values assigned by *scikit-learn* library. The link to the GitHub folder with the developed codes is given in the appendix.

2.4 Results analysis

In the literature coefficient of determination (R^2) seems to be the go-to statistical parameter to evaluate the performance of the RFR models.⁶¹⁻⁶³ However, an improved version of R^2 known as adjusted coefficient of determination (adj. R^2) was used in this study to evaluate the developed



models since it takes the number of datapoints and the number of input features into consideration during evaluation.⁵ The mathematical equation of R^2 is shown in Eq. 7 and Eq. 8 shows the mathematical equation of adj. R^2 . The closer the adj. R^2 to 1, the higher the performance of the model.

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}. \quad (7)$$

$$Adj. R^2 = 1 - \frac{(1-R^2)(n-1)}{n-m-1}. \quad (8)$$

In Eq. 7 and Eq. 8 y_i is the actual value, \hat{y} is the predicted value, \bar{y} is the mean value of the distribution, n is the number of datapoints and m is the number of input features.

Apart from the adj. R^2 , models were evaluated using RMSE as well. The mathematical equation of RMSE is shown in Eq. 9.

$$RMSE = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n}}. \quad (9)$$

developed RFR models were further evaluated based on their runtime to study how tuning considered hyperparameters affects computational times. Train-test difference was also used to further analyze the models. Train-test difference is an indication of the generalizability of an ML model, and it gives an idea about the variance of the model. The lower the train-test difference, the higher the generalizability of the model.^{64,65}

3. Results and discussion

The adj. R^2 values obtained using approach 1 of the methodology are tabulated in Table 3. When only $n_{estimators}$ increased from 1 to 10 in intervals of 1 (keeping all the other hyperparameters at their default values), the model yielded a training adj. R^2 of 0.9650, validation adj. R^2 of 0.8188



and testing adj. R^2 of 0.8024. The training score is higher than the validation (cross-validation) and testing scores as expected since the model is fitted (trained) to the training set and this pattern was observed in all the models developed during the study. In approach 1, when the upper limit of $n_{estimators}$ value was increased from 10 to 100, training, validation and testing scores showed a significant increase. Training score had an increase of 1.14%, validation score had an increase of 2.19%. The test score had an increase of 2.22%. This rise in performance can be clearly seen in Fig. 7 where the adj. R^2 values of the testing models were plotted for each approach.

Table 3 Coefficient of determination, train-test difference and computational times of the models obtained in approach 1.

| Model No. | $n_{estimator}$ r Change | $n_{estimator}$ | adj. R^2 | | | Computation Time (sec) |
|-----------|--------------------------|-----------------|------------|-----------|---------|------------------------|
| | | | Training | Validatio | Testing | |
| M11 | 1:10:1 | 8 | 0.9650 | 0.8188 | 0.8024 | 0.81 |
| M12 | 1:100:1 | 51 | 0.9760 | 0.8367 | 0.8202 | 70.25 |
| M13 | 1:100:10 | 51 | 0.9760 | 0.8367 | 0.8202 | 6.88 |
| M14 | 1:1000:1 | 51 | 0.9760 | 0.8367 | 0.8202 | 6932.55 |
| M15 | 1:1000:10 | 51 | 0.9760 | 0.8367 | 0.8202 | 707.56 |
| M16 | 1:1000:100 | 801 | 0.9799 | 0.8352 | 0.8218 | 65.73 |



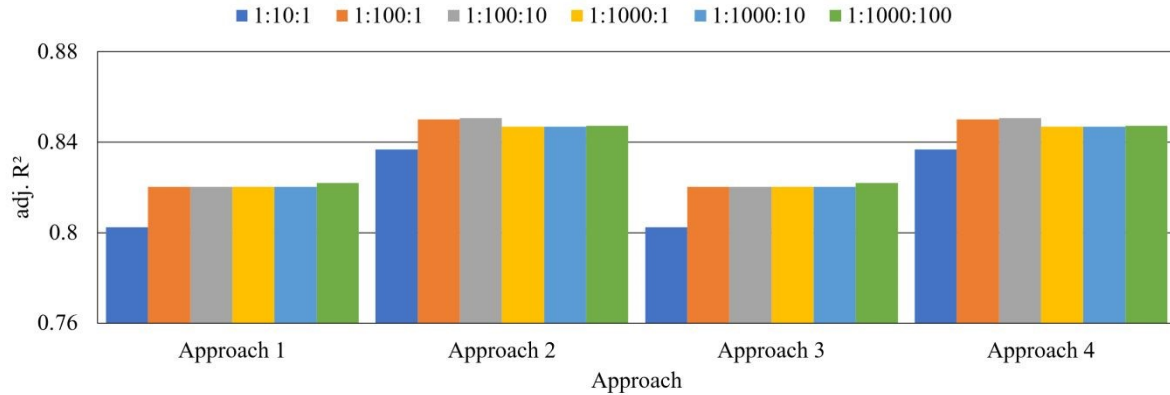


Fig. 7 Coefficient of determination values of each approach for different $n_{estimators}$ change.

Interestingly, when the upper limit of the $n_{estimators}$ range was pushed beyond 100, performance of the model did not show any noticeable increase in all training validation and testing adj. R^2 values. In fact, when $n_{estimators}$ changed from 1 to 100 in intervals 1 and 10 (models M12 and M13) and $n_{estimators}$ changed from 1 to 1000 in intervals 1 and 10 (models M14 and M15) models showed the same performance, i.e. a training score of 0.9760, validation score of 0.8367 and a testing score of 0.8202. However, when $n_{estimators}$ were changed from 1 to 1000 in intervals of 100, training and testing scores of the model M16 showed a slight increase of performance yielding an adj. R^2 of 0.9799 and 0.8218 respectively. However, validation score showed a slight decrease which was negligible.

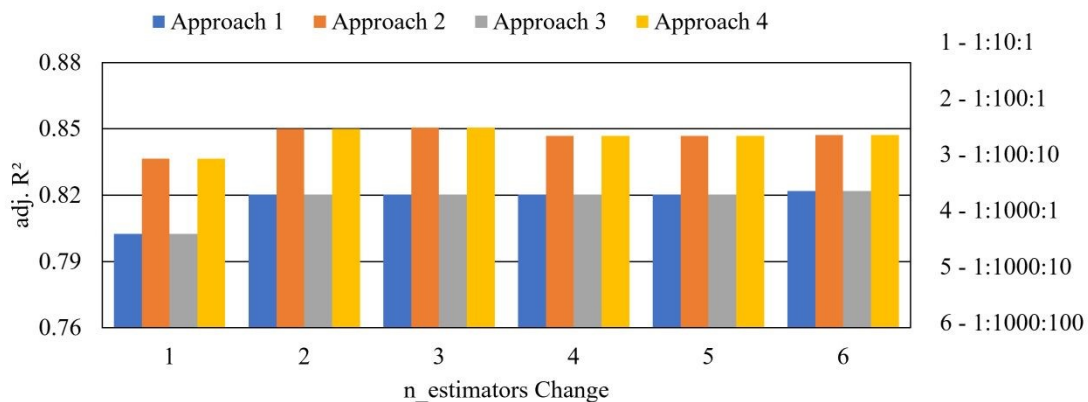


Fig. 8 Coefficient of determination values of each $n_{estimators}$ change for different approaches.



Further, the highest computational time of 6932.55 seconds was shown by the model M14 where $n_estimators$ was changed from 1 to 1000 in increments of 1. The results from approach 1 showed that after a certain $n_estimators$ value models' performances increased drastically and the performance was maintained at a constant value over a certain $n_estimators$ range showing that the performance of the RFR when $n_estimators$ were tuned was efficient within a certain range. Since the range and interval at which the $n_estimators$ values are tuned affects the computational time an effective range and an interval for $n_estimators$ to be decided upon taking computational time into account.

In approach 2, $max_features$ were also tuned along with $n_estimators$. Results obtained using approach 2 of the methodology are tabulated in Table 4. As observed in approach 1, a clear spike in training, validation and testing adj. R^2 values was observed when $n_estimators$ ' upper limit was increased from 10 to 100. Training score had an increase of 1.36%, validation score had an increase of 1.92%. The test score had an increase of 1.60%. This clear jump in performance is noticeable in Fig. 7. Interestingly, the performances of the models developed in approach 2 were significantly higher than the performance of the corresponding " $n_estimators$ change" in approach 1. This is quite visible in Fig. 8 as well. Further, going from approach 1 to 2, the average validation score increased by 2.24% and the testing score increased by 3.52% which was significant. This increase of adj. R^2 values is an indication that tuning $max_features$ have a major impact on the predicting porosity using RFR. Model M21 where $n_estimators$ were changed from 1 to 10 in intervals of 1 and $max_features$ were changed from 0.1 to 1 in intervals of 0.1 showed the least performance with a training score of 0.9672, validation score of 0.8381 and a testing score of 0.8366. On the other hand, model M23 showed the highest testing performance with an adj. R^2 of 0.8505 where $n_estimators$ were changed from 1 to 100 in intervals of 10 and $max_features$ were changed from



0.1 to 1 in intervals of 0.1. The model M23 yielded its best test model when $n_estimators$ were 81 and $max_features$ were 0.5. It has to be noted that even though the model M23 had the highest testing score, training and validation scores were not the best out of all the models developed in approach 2. The highest training score of 0.9823 was shown by models M24, M25 and M26. Highest validation scores were shown by models M24 and M25. However, it is more meaningful to select model M23 as the best performing model since the testing set represents an independent dataset which had never been seen by the model before.

Table 4 Coefficient of determination, train-test difference and computational times of the models obtained in approach 2.

| Model No. | n_estimato r Change | n_estimator | max_feature s | adj. R ² | | | Computation Time (sec) |
|-----------|------------------------|-------------|------------------|---------------------|----------------|---------|---------------------------|
| | | | | Training | Validatio n | Testing | |
| M21 | 1:10:1 | 9 | 0.1 | 0.9672 | 0.8381 | 0.8366 | 3.69 |
| M22 | 1:100:1 | 79 | 0.5 | 0.9804 | 0.8542 | 0.8500 | 326.56 |
| M23 | 1:100:10 | 81 | 0.5 | 0.9806 | 0.8541 | 0.8505 | 30.20 |
| M24 | 1:1000:1 | 520 | 0.5 | 0.9823 | 0.8556 | 0.8467 | 32620.39 |
| M25 | 1:1000:10 | 521 | 0.5 | 0.9823 | 0.8556 | 0.8467 | 3045.27 |
| M26 | 1:1000:100 | 801 | 0.5 | 0.9823 | 0.8554 | 0.8471 | 284.29 |

The anomaly in the validation score observed when the $n_esitimators$ were changed from 1 to 1000 in intervals of 100 in approach 1 was observable in approach 2 as well. The difference in train-test score gives an idea about the generalizability of the model. The lesser the train-test difference, the higher the generalizability of the model. Overall train-test difference of approach 2 was noticeably less than that of approach 1. The average train-test difference decreased by 15.51%



going from approach 1 to 2. This showed that the generalizability of the models improved when *max_features* was introduced to the hyperparameter space. Similar to that of approach 1 highest runtime was shown when the *n_estimators* were changed from 1 to 1000 in increments of 1.

In approach 3 *n_estimators* was investigated with the alteration of *min_samples_leaf* and the results obtained are tabulated in Table 5. Noticeably, all the performance results obtained for all the RFR models except the runtimes were the same as that of approach 1 as seen in Fig. 7 and Fig. 8. The reason for this was the optimum value selected by grid search optimization of *min_samples_leaf* was same as the default value assigned by *scikit-learn* library for the RFR algorithm, hence the best testing adj. R^2 was shown by model M34 when *n_estimators* were changed from 1 to 1000 in intervals of 100. Computational times were clearly higher than those obtained in approach 1 since models developed in approach 3 had a larger hyperparameter space compared to approach 1.

Table 5 Coefficient of determination, train-test difference and computational times of the models obtained in approach 3.

| Model No. | n_estimator Change | n_estimators | min_samples_leaf | adj. R^2 | | | Computation Time (sec) |
|-----------|--------------------|--------------|------------------|------------|------------|---------|------------------------|
| | | | | Training | Validation | Testing | |
| M31 | 1:10:1 | 8 | 1 | 0.9650 | 0.8188 | 0.8024 | 7.79 |
| M32 | 1:100:1 | 51 | 1 | 0.9760 | 0.8367 | 0.8202 | 674.81 |
| M33 | 1:100:10 | 51 | 1 | 0.9760 | 0.8367 | 0.8202 | 64.96 |
| M34 | 1:1000:1 | 51 | 1 | 0.9760 | 0.8367 | 0.8202 | 70039.55 |
| M35 | 1:1000:10 | 51 | 1 | 0.9760 | 0.8367 | 0.8202 | 6525.18 |
| M36 | 1:1000:100 | 801 | 1 | 0.9799 | 0.8352 | 0.8218 | 606.28 |



In approach 4, $n_estimators$ was changed along with both $max_features$ and $min_samples_leaf$. Results in Table 6 showed that the performances of the models were the same as that of approach 2. A similar phenomenon caused this performance similarity as observed between approach 1 and approach 3. In this case, $min_samples_leaf$ always selected the default value during the tuning process and the $max_features$ selected for the optimum model was similar to that of approach 2. Approach 4 consumed the longest computational time since 3 hyperparameters had to be tuned simultaneously. The highest runtime consumed for all models was recorded in this approach by the model M44, which was 82832.02 seconds. In approach 4 as observed in approach 2 as well, there was a test model performance increase of 1.60% when the upper limit of $n_estimators$ was increased from 10 to 100. When the upper limit was increased from 100 to 1000 there was a test model performance drop of around 0.4%.

Table 6 Coefficient of determination, train-test difference and computational times of the models obtained in approach 4.

| Model No. | n_estimator Change | n_estimator | max_features | min_samples_leaf | adj. R ² | | | Computation Time (sec) |
|-----------|--------------------|-------------|--------------|------------------|---------------------|----------------|---------|------------------------|
| | | | | | Trainin g | Validatio n | Testing | |
| M41 | 1:10:1 | 9 | 0.1 | 1 | 0.9672 | 0.8381 | 0.8366 | 56.22 |
| M42 | 1:100:1 | 79 | 0.5 | 1 | 0.9804 | 0.8542 | 0.8500 | 4242.86 |
| M43 | 1:100:10 | 81 | 0.5 | 1 | 0.9806 | 0.8541 | 0.8505 | 425.65 |
| M44 | 1:1000:1 | 520 | 0.5 | 1 | 0.9823 | 0.8556 | 0.8467 | 82832.02 |
| M45 | 1:1000:10 | 521 | 0.5 | 1 | 0.9823 | 0.8556 | 0.8467 | 51444.27 |
| M46 | 1:1000:100 | 801 | 0.5 | 1 | 0.9823 | 0.8554 | 0.8471 | 3796.99 |



Table 7 RMSE of training and testing models of approaches 1, 2, 3 and 4.

| RMSE | | | | | | | | | | | |
|------------|-----------|----------|------------|------------|----------|------------|-----------|---------|------------|----------|----------|
| Approach 1 | | | Approach 2 | | | Approach 3 | | | Approach 4 | | |
| Model No. | Trainin g | Testin g | Model No. | Training g | Testin g | Model No. | Trainin g | Testing | Model No. | Training | Testin g |
| M11 | 1.2894 | 2.9967 | M21 | 1.2516 | 2.7218 | M31 | 1.2894 | 2.9967 | M41 | 1.2516 | 2.7218 |
| M12 | 1.0817 | 2.8499 | M22 | 0.9835 | 2.5917 | M32 | 1.0817 | 2.8499 | M42 | 0.9835 | 2.5917 |
| M13 | 1.0817 | 2.8499 | M23 | 0.9798 | 2.5875 | M33 | 1.0817 | 2.8499 | M43 | 0.9798 | 2.5880 |
| M14 | 1.0817 | 2.8499 | M24 | 0.9399 | 2.6190 | M34 | 1.0817 | 2.8499 | M44 | 0.9399 | 2.6190 |
| M15 | 1.0817 | 2.8499 | M25 | 0.9396 | 2.6187 | M35 | 1.0817 | 2.8499 | M45 | 0.9396 | 2.6187 |
| M16 | 0.9988 | 2.8312 | M26 | 0.9396 | 2.6148 | M36 | 0.9988 | 2.8312 | M46 | 0.9396 | 2.6148 |

Table 7 shows the RMSE values of approaches 1, 2, 3 and 4. While adj. R^2 values give an idea about the correlation between the actual porosities and the predicted porosities, the RMSE values give an idea about the difference (or the error) between the two. Therefore, RMSE is also an important parameter in ML model performance evaluation. The pattern in which RMSE values fluctuated in the 4 approaches was similar to that of adj. R^2 . The least RMSE's were shown by models M16 with a training model RMSE of 0.9988 and a testing model RMSE of 2.8312. The results improvement when *max_features* were introduced to the hyperparameter space was also evident based on the RMSE values obtained in approach 2. There was a clear decrease in RMSE values in both training and testing models in approach 2 and 4 where *max_features* was tuned.



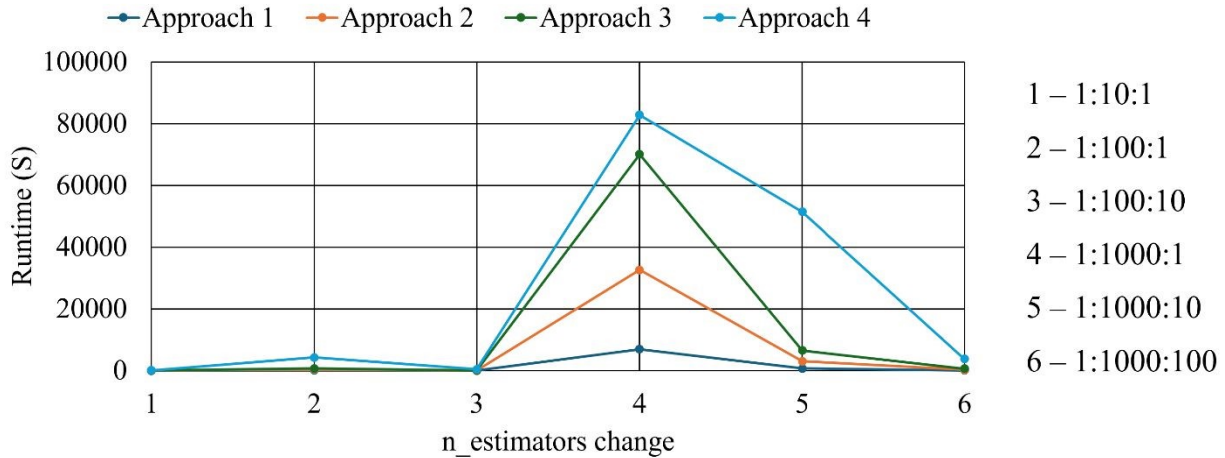


Fig. 9 Runtime of the models of each $n_estimators$ change for different approaches.

Runtime and grid search combinations had a positive relationship, i.e., when the number of combinations in the grid search space was the largest, the runtime of the model was the highest and vice versa. Further, it was observed that, from approach 1 to approach 3, the increase of computational times was roughly proportional to each other as seen in Fig. 9. However, in approach 4, where $n_estimators$ was changed along with the tuning of $max_features$ and $min_samples_leaf$, an anomaly was observed when $n_estimators$ were changed from 1 to 1000 in intervals of 10.



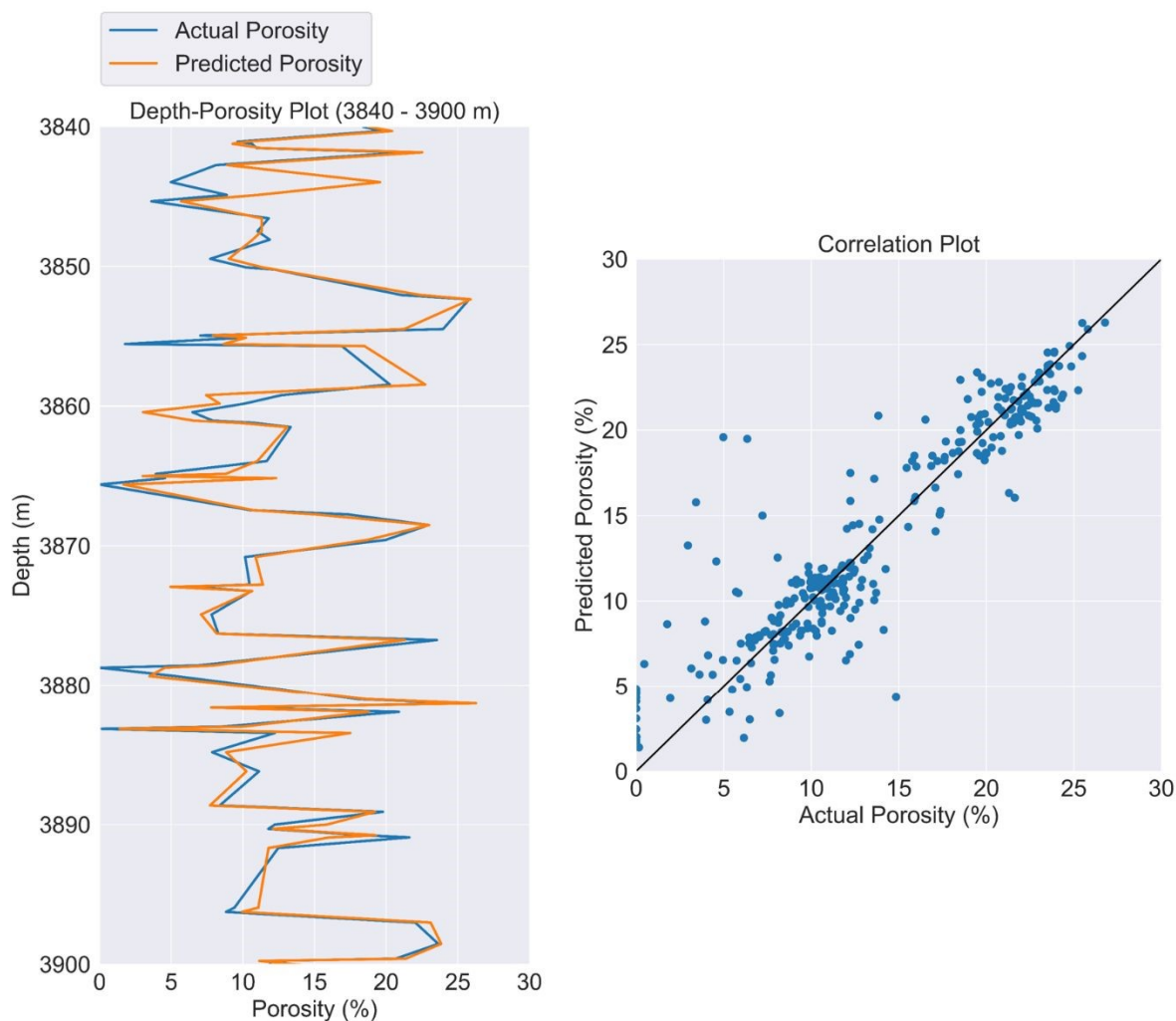


Fig 10. Depth-porosity and correlation plots obtained from the predictions of the best performing RFR testing model.

Even though the primary objective of the study was to investigate the influence of $n_estimators$ along with $max_features$ and $min_samples_leaf$ on the performance of RFR, having an overall picture on the variation of actual and predicted porosity and their relationship is important to understand the model's applicability in porosity prediction. To achieve this, depth-porosity graphs and correlation plots were plotted. Fig. 10 shows one such depth-porosity graph and a correlation plot developed for the best performing RFR test model (model M23) of the study. Depth-porosity plot indicated that most of the time the predicted porosity followed the pattern of the actual



porosity. The correlation plot showed that the majority of the points scattered around the perfect correlation line which is an indication of a high correlation between the actual values and the predicted values.

4. Conclusions

The aim of this study was to examine the effects of tuning the number of decision trees in the forest ($n_estimators$) in random forest regression (RFR) for predicting porosity within the Volve oil field in North Sea. Additionally, the study investigated the influence of $n_estimators$ when tuned with two others commonly used hyperparameters, namely, number of features to consider when looking for the best split ($max_features$) and minimum number of samples required to be at a leaf node ($min_samples_leaf$). The hyperparameters were tuned using grid search optimization integrating 5-fold cross validation and model performances were evaluated based on adj. R^2 , RMSE and computational times.

- Overall, based on both performance and computational time, the RFR model with $n_estimators$ at 81 and $max_features$ at 2 (while keeping all the other hyperparameters at their default values), which was developed in approach 2, produced the most effective model for predicting porosity of Volve oil field in North Sea with a testing model adj. R^2 of 0.8505, a testing model RMSE of 2.5875 and a computational time of 30.2 seconds.
- There was a notable increase in performance when the upper limit of the $n_estimators$ increased from 10 to 100. On the other hand, the performance of the models did not increase significantly when the upper limit of $n_estimators$ increased from 100 to 1000. This phenomenon indicated that identifying an effective $n_estimators$ range which is not too low (which will make the performance significantly low) and not too high (which will increase the computational time) is important to produce an efficient RFR model during porosity prediction.



- Further, range 1 to 100 changed in intervals of 10 can be suggested for $n_estimators$ when developing an RFR model to predict porosity of Volve oil field, since those models showed higher performances and lower computational times in all four approaches. When $n_estimators$ range 1 to 100 was changed in intervals of 10 it always yielded a high adj. R^2 value (in approach 2 and 4 it yielded the highest testing model adj. R^2 value) for the model and consumed the second least computational time.
- When $n_estimators$ tuned along with $max_features$ in approach 2 the results improved drastically compared to approach 1 where only $n_estimators$ were tuned. There was an average validation score increase of 2.24% and the testing score increase of 3.52% going from approach 1 to 2. This improvement of the scores (adj. R^2) showed that $max_features$ have a significant influence on the RFR model's performance.
- Moreover, it was observed that computational time was largely affected by the number of hyperparameters altered, their range and interval. Out of all the approaches, the highest computational time was consumed when $n_estimators$ were tuned from 1 to 1000 in intervals of 1 along with $max_features$ and $min_samples_leaf$.

Based on the results, only by adjusting $n_estimators$ and $max_features$ an RFR model can be developed with a robust prediction power to estimate porosity in Volve oil field.

Recommendations

This study focused on three hyperparameters, namely, $n_estimators$, $max_features$ and $min_samples_leaf$. Apart from these hyperparameters $min_samples_split$ and max_depth are also widely used in the literature during hyperparameter optimization in RFR. Therefore, for future



studies the behaviour of *min_samples_split* and *max_depth* along with *n_estimators* are recommended to be investigated.

Author contributions

Kushan Sandunil: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Validation, Visualization, Writing – original draft. Ziad Bennour: Conceptualization, Methodology, Supervision, Funding acquisition, Resources, Writing - review & editing. Hisham Ben Mahmud: Supervision, Funding acquisition, Writing - review & editing. Ausama Giwelli: Supervision, Writing - review & editing.

Declaration of competing interest

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

Data availability statement

Data for this article, including codes and graphs are available at GitHub at <https://github.com/kwkushan/effects-of-tuning-decision-trees-in-random-forest-regression-on-predicting-porosity-kushan-sandunil->. A description about the available files in GitHub repository can be found in the appendix.

Acknowledgement

Authors would like to thank Curtin University Malaysia and Curtin Malaysia Postgraduate Research Scholarship (CMPRS) for hosting and allocating the research grant for the study. Further, a special thanks would be given to Equinor and the Volve license partners for making the Volve field dataset available for scientific research (<https://discovervolve.com/citation-non-commerciality-clause/>).



Appendix

The authors would like to share open repository folder containing the codes and resources for this study on [GitHub](#) and extend an invitation to collaborate through Open Knowledge sharing. In the folder, 4 codes are provided; Code_1 was developed by only tuning $n_estimators$. Code_2 was developed by tuning $n_estimators$ along with $max_features$. Code_3 was developed by tuning $n_estimators$ along with $min_samples_leaf$ and the Code_4 was developed by tuning all three hyperparameters, i.e., $n_estimators$, $max_features$ and $min_samples_leaf$. The GitHub repository with the developed codes in the study can be accessed via this [link](#).

Nomenclature

| | |
|------------|---|
| AI | Artificial Intelligence |
| ML | Machine Learning |
| RFR | Random Forest Regression |
| ANN | Artificial Neural Network |
| SVR | Support Vector Regression |
| LASSO | Least Absolute Shrinkage and Selection Operator |
| BMA | Bayesian Model Averaging |
| GSO | Grid Search Optimization |
| RMSE | Root Mean Squared Error |
| R^2 | Coefficient of determination |
| adj. R^2 | Adjusted coefficient of determination |
| RES | Resistivity log |
| NPHI | Neutron porosity log |
| GR | Gamma ray log |
| PHIF | Total porosity |
| PHID | Porosity from density log |



| | |
|-----------------------|--|
| $n_estimators$ | Number of trees in the forest |
| $max_features$ | Number of features considered for the best split |
| $min_samples_leaf$ | Minimum number of samples required to be at a leaf node |
| max_depth | Maximum depth of the tree |
| $max_samples$ | Fraction of the original dataset assigned to any individual tree |
| $min_samples_split$ | Minimum number of samples required to split an internal node |
| max_leaf_nodes | Maximum leaf nodes to restrict the growth of the tree |
| A | A regression coefficient |
| B | A regression coefficient |
| ρ_{ma} | Matrix density |
| ρ_b | Measured bulk density |
| ρ_{fl} | Pore fluid density |
| n | Number of datapoints |
| m | Number of input features |
| X | Independent and identically distributed random vector |
| θ_r | Independent and identically distributed random vector |
| x | Observed input vector associated with vector X |
| Y | A vector with numerical outcomes |
| y_i | Actual Value |
| \hat{y} | Predicted Value |
| \bar{y} | Mean Value of the Distribution |

References

1. Kavuri C, Kokjohn SL. Exploring the potential of machine learning in reducing the computational time/expense and improving the reliability of engine optimization studies. *International Journal of Engine Research*. 2020 Sep;21(7):1251-70.
2. Zhan N, Kitchin JR. Uncertainty quantification in machine learning and nonlinear least squares regression models. *AIChE Journal*. 2022 Jun;68(6):e17516.



3. Zhang X, Tian Y, Chen L, Hu X, Zhou Z. Machine learning: a new paradigm in computational electrocatalysis. *The Journal of Physical Chemistry Letters*. 2022 Aug 18;13(34):7920-30.
4. Turing AM. *Computing machinery and intelligence*. Springer Netherlands; 2009.
5. Al-Mudhafar WJ. Bayesian and LASSO regressions for comparative permeability modeling of sandstone reservoirs. *Natural Resources Research*. 2019 Jan 15;28(1):47-62.
6. Ojukwu C, Smith K, Kadkhodayan N, Leung M, Baldwin K. Reservoir Characterization, Machine Learning and Big Data—An Offshore California Case Study. *InSPE Nigeria Annual International Conference and Exhibition 2020* Aug 11 (p. D013S002R005). SPE.
7. Silva AA, Tavares MW, Carrasquilla A, Misságia R, Ceia M. Petrofacies classification using machine learning algorithms. *Geophysics*. 2020 Jul 1;85(4):WA101-13.
8. Amiri M, Ghiasi-Freez J, Golkar B, Hatampour A. Improving water saturation estimation in a tight shaly sandstone reservoir using artificial neural network optimized by imperialist competitive algorithm—A case study. *Journal of Petroleum Science and Engineering*. 2015 Mar 1;127:347-58.
9. Elkatatny S, Mahmoud M, Tariq Z, Abdulraheem A. New insights into the prediction of heterogeneous carbonate reservoir permeability from well logs using artificial intelligence network. *Neural Computing and Applications*. 2018 Nov;30:2673-83.
10. Akande KO, Owolabi TO, Olatunji SO, AbdulRaheem A. A hybrid particle swarm optimization and support vector regression model for modelling permeability prediction of hydrocarbon reservoir. *Journal of Petroleum Science and Engineering*. 2017 Feb 1;150:43-53.
11. Baziar S, Shahripour HB, Tadayoni M, Nabi-Bidhendi M. Prediction of water saturation in a tight gas sandstone reservoir by using four intelligent methods: a comparative study. *Neural Computing and Applications*. 2018 Aug;30:1171-85.
12. Anifowose F, Abdulraheem A, Al-Shuhail A. A parametric study of machine learning techniques in petroleum reservoir permeability prediction by integrating seismic attributes and wireline data. *Journal of Petroleum Science and Engineering*. 2019 May 1;176:762-74.
13. Kamali MZ, Davoodi S, Ghorbani H, Wood DA, Mohamadian N, Lajmorak S, Rukavishnikov VS, Taherizade F, Band SS. Permeability prediction of heterogeneous carbonate gas condensate reservoirs applying group method of data handling. *Marine and Petroleum Geology*. 2022 May 1;139:105597.



14. Al-Mudhafar W. Integrating bayesian model averaging for uncertainty reduction in permeability modeling. Inoffshore technology conference 2015 May 4 (pp. OTC-25646). OTC.
15. Wang G, Ju Y, Li C, Carr TR, Cheng G. Application of artificial intelligence on black shale lithofacies prediction in Marcellus Shale, Appalachian Basin. InUnconventional Resources Technology Conference, Denver, Colorado, 25-27 August 2014 2014 Aug 27 (pp. 1970-1980). Society of Exploration Geophysicists, American Association of Petroleum Geologists, Society of Petroleum Engineers.
16. Al-Mudhafar WJ. Integrating well log interpretations for lithofacies classification and permeability modeling through advanced machine learning algorithms. *Journal of Petroleum Exploration and Production Technology*. 2017 Dec;7(4):1023-33.
17. J Al-Mudhafar W. Integrating lithofacies and well logging data into smooth generalized additive model for improved permeability estimation: Zubair formation, South Rumaila oil field. *Marine Geophysical Research*. 2019 Sep;40:315-32.
18. Kim J. Lithofacies classification integrating conventional approaches and machine learning technique. *Journal of Natural Gas Science and Engineering*. 2022 Apr 1;100:104500.
19. Na'imi SR, Shadizadeh SR, Riahi MA, Mirzakhani M. Estimation of reservoir porosity and water saturation based on seismic attributes using support vector regression approach. *Journal of Applied Geophysics*. 2014 Aug 1;107:93-101.
20. Al-AbdulJabbar A, Al-Azani K, Elkatatny S. Estimation of reservoir porosity from drilling parameters using artificial neural networks. *Petrophysics*. 2020 Jun 1;61(03):318-30.
21. Chen W, Yang L, Zha B, Zhang M, Chen Y. Deep learning reservoir porosity prediction based on multilayer long short-term memory network. *Geophysics*. 2020 Jul 1;85(4):WA213-25.
22. Anifowose FA. Ensemble machine learning: the latest development in computational intelligence for petroleum reservoir characterization. InSPE Kingdom of Saudi Arabia Annual Technical Symposium and Exhibition 2013 May 19 (pp. SPE-168111). SPE.
23. Subasi A, El-Amin MF, Darwich T, Dossary M. Permeability prediction of petroleum reservoirs using stochastic gradient boosting regression. *Journal of Ambient Intelligence and Humanized Computing*. 2022 Apr:1-0.



24. Wang J, Yan W, Wan Z, Wang Y, Lv J, Zhou A. Prediction of permeability using random forest and genetic algorithm model. *Computer Modeling in Engineering & Sciences*. 2020 Dec 10;125(3):1135-57.
25. Otchere DA, Ganat TO, Gholami R, Lawal M. A novel custom ensemble learning model for an improved reservoir permeability and water saturation prediction. *Journal of Natural Gas Science and Engineering*. 2021 Jul 1;91:103962.
26. Zhang Z, Cai Z. Permeability prediction of carbonate rocks based on digital image analysis and rock typing using random forest algorithm. *Energy & Fuels*. 2021 Jun 24;35(14):11271-84.
27. Lee TH, Ullah A, Wang R. Bootstrap aggregating and random forest. *Macroeconomic forecasting in the era of big data: Theory and practice*. 2020:389-429.
28. Maher MM, Sakr S. Smartml: A meta learning-based framework for automated selection and hyperparameter tuning for machine learning algorithms. In *EDBT: 22nd International conference on extending database technology* 2019 Mar 26.
29. Yang L, Shami A. On hyperparameter optimization of machine learning algorithms: Theory and practice. *Neurocomputing*. 2020 Nov 20;415:295-316.
30. Isabona J, Imoize AL, Kim Y. Machine learning-based boosted regression ensemble combined with hyperparameter tuning for optimal adaptive learning. *Sensors*. 2022 May 16;22(10):3776.
31. Zou C, Zhao L, Xu M, Chen Y, Geng J. Porosity prediction with uncertainty quantification from multiple seismic attributes using random forest. *Journal of Geophysical Research: Solid Earth*. 2021 Jul;126(7):e2021JB021826
32. Rezaee R, Ekundayo J. Permeability prediction using machine learning methods for the CO₂ injectivity of the precipice sandstone in Surat Basin, Australia. *Energies*. 2022 Mar 11;15(6):2053.
33. García S, Luengo J, Herrera F. *Data preprocessing in data mining*. Cham, Switzerland: Springer International Publishing; 2015 May.
34. Gudivada V, Apon A, Ding J. Data quality considerations for big data and machine learning: Going beyond data cleaning and transformations. *International Journal on Advances in Software*. 2017 Jul;10(1):1-20.
35. Maharana K, Mondal S, Nemade B. A review: Data pre-processing and data augmentation techniques. *Global Transitions Proceedings*. 2022 Jun 1;3(1):91-9.



36. Al Ghaithi A, Prasad M. Machine learning with artificial neural networks for shear log predictions in the Volve field Norwegian North Sea. InSEG Technical Program Expanded Abstracts 2020 2020 Sep 30 (pp. 450-454). Society of Exploration Geophysicists.
37. Ng CS, Ghahfarokhi AJ, Amar MN. Well production forecast in Volve field: Application of rigorous machine learning techniques and metaheuristic algorithm. *Journal of Petroleum Science and Engineering*. 2022 Jan 1;208:109468.
38. Nikitin NO, Revin I, Hvatov A, Vychuzhanin P, Kalyuzhnaya AV. Hybrid and automated machine learning approaches for oil fields development: The case study of Volve field, North Sea. *Computers & Geosciences*. 2022 Apr 1;161:105061.
39. Mapchart. World map: simple [Internet]. 2024 [cited 2024 Jul 22]. Available from: <https://www.mapchart.net/world.html>
40. Sen S, Ganguli SS. Estimation of pore pressure and fracture gradient in Volve field, Norwegian North Sea. InSPE Oil and Gas India Conference and Exhibition. 2019 Apr 8 (p. D022S027R002). SPE.
41. Statoil. 15/9-19A Well Composite Log, Sleipner, Theta Vest Prospect Structure [Internet]. 1998 [cited 2023 Mar 1]. Available from: <https://discovervolve.com/citation-non-commerciality-clause/>
42. Ilyas IF, Chu X. Data cleaning. Morgan & Claypool; 2019 Jun 18.
43. Jain A, Patel H, Nagalapatti L, Gupta N, Mehta S, Guttula S, Mujumdar S, Afzal S, Sharma Mittal R, Munigala V. Overview and importance of data quality for machine learning tasks. InProceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining 2020 Aug 23 (pp. 3561-3562).
44. Rawat S, Rawat A, Kumar D, Sabitha AS. Application of machine learning and data visualization techniques for decision support in the insurance sector. *International Journal of Information Management Data Insights*. 2021 Nov 1;1(2):100012.
45. Ahamad MM, Aktar S, Rashed-Al-Mahfuz M, Uddin S, Liò P, Xu H, Summers MA, Quinn JM, Moni MA. A machine learning model to identify early stage symptoms of SARS-Cov-2 infected patients. *Expert systems with applications*. 2020 Dec 1;160:113661.
46. Sarker IH, Abushark YB, Alsolami F, Khan AI. Intrudtree: a machine learning based cyber security intrusion detection model. *Symmetry*. 2020 May 6;12(5):754.



47. Feizi H, Apaydin H, Sattari MT, Colak MS, Sibtain M. Improving reservoir inflow prediction via rolling window and deep learning-based multi-model approach: case study from Ermenek Dam, Turkey. *Stochastic Environmental Research and Risk Assessment*. 2022 Oct;36(10):3149-69.
48. Salazar JJ, Garland L, Ochoa J, Pyrcz MJ. Fair train-test split in machine learning: Mitigating spatial autocorrelation for improved prediction accuracy. *Journal of Petroleum Science and Engineering*. 2022 Feb 1;209:109885.
49. Mask GM, Wu X. Deriving New Type Curves through Machine Learning in the Wolfcamp Formation. In *SPE Reservoir Characterisation and Simulation Conference and Exhibition*. 2023 Jan 24 (p. D011S001R007). SPE.
50. Emmanuel T, Maupong T, Mpoeleng D, Semong T, Mphago B, Tabona O. A survey on missing data in machine learning. *Journal of Big data*. 2021 Dec;8:1-37.
51. Seliem MM. Handling Outlier data as missing values by imputation methods: application of machine learning algorithms. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*. 2022 Jan 30;13(1):273-86.
52. Garcia-Carretero R, Holgado-Cuadrado R, Barquero-Pérez Ó. Assessment of classification models and relevant features on nonalcoholic steatohepatitis using random forest. *Entropy*. 2021 Jun 17;23(6):763.
53. Suherman IC, Sarno R. Implementation of random forest regression for COCOMO II effort estimation. In *2020 international seminar on application for technology of information and communication (iSemantic) 2020 Sep 19 (pp. 476-481)*. IEEE.
54. Yilmazer S, Kocaman S. A mass appraisal assessment study using machine learning based on multiple regression and random forest. *Land use policy*. 2020 Dec 1;99:104889.
55. Segal MR. Machine learning benchmarks and random forest regression.
56. Abbaszadeh M, Soltani-Mohammadi S, Ahmed AN. Optimization of support vector machine parameters in modeling of Iju deposit mineralization and alteration zones using particle swarm optimization algorithm and grid search method. *Computers & Geosciences*. 2022 Aug 1;165:105140.
57. Abbas MA, Al-Mudhafar WJ, Wood DA. Improving permeability prediction in carbonate reservoirs through gradient boosting hyperparameter tuning. *Earth Science Informatics*. 2023 Dec;16(4):3417-32.



58. Sandunil K, Bennour Z, Ben Mahmud H, Giwelli A. Effects of Tuning Hyperparameters in Random Forest Regression on Reservoir's Porosity Prediction. Case Study: Volve Oil Field, North Sea. In ARMA US Rock Mechanics/Geomechanics Symposium 2023 Jun 25 (pp. ARMA-2023). ARMA.
59. Al-Mudhafar WJ. Incorporation of bootstrapping and cross-validation for efficient multivariate facies and petrophysical modeling. In SPE Rocky Mountain Petroleum Technology Conference/Low-Permeability Reservoirs Symposium 2016 May 5 (pp. SPE-180277). SPE.
60. Rahimi M, Riahi MA. Reservoir facies classification based on random forest and geostatistics methods in an offshore oilfield. *Journal of Applied Geophysics*. 2022 Jun 1;201:104640.
61. Mahmoud AA, Elkatatny S, Chen W, Abdurraheem A. Estimation of oil recovery factor for water drive sandy reservoirs through applications of artificial intelligence. *Energies*. 2019 Sep 25;12(19):3671.
62. Al Khalifah H, Glover PW, Lorinczi P. Permeability prediction and diagenesis in tight carbonates using machine learning techniques. *Marine and Petroleum Geology*. 2020 Feb 1;112:104096.
63. Ridwan WM, Sapitang M, Aziz A, Kushiari KF, Ahmed AN, El-Shafie A. Rainfall forecasting model using machine learning methods: Case study Terengganu, Malaysia. *Ain Shams Engineering Journal*. 2021 Jun 1;12(2):1651-63.
64. Wang H, Lei Z, Zhang X, Zhou B, Peng J. Machine learning basics. *Deep learning*. 2016:98-164.
65. Mehta P, Bukov M, Wang CH, Day AG, Richardson C, Fisher CK, Schwab DJ. A high-bias, low-variance introduction to machine learning for physicists. *Physics reports*. 2019 May 30;810:1-24.



Data availability statement

Data for this article, including codes and graphs are available at GitHub at <https://github.com/kwkushan/effects-of-tuning-decision-trees-in-random-forest-regression-on-predicting-porosity-kushan-sandunil->. A description about the available files in GitHub repository can be found in the appendix.

