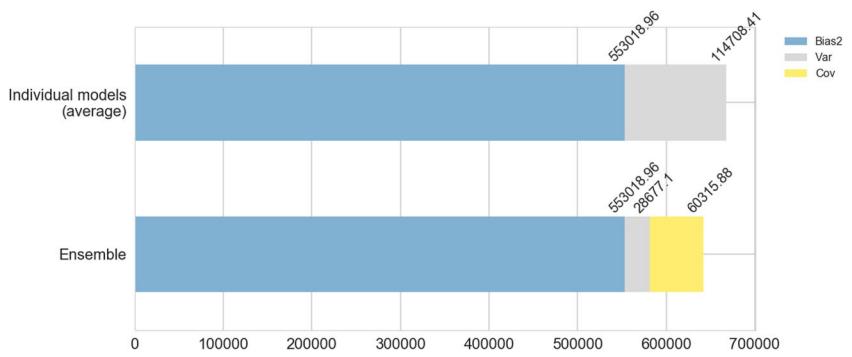


Fig. 1 Squared error bias (squared bias—Bias2), variance (Var), and covariance (Cov) decomposition for an heterogeneous ensemble of surrogates—Rosenbrock 6D test function



surrogate predictions based on performance (error) measures that are either static/global (Goel et al. 2007; Acar and Rais-Rohani 2009), dynamic/local (Zerpa et al. 2005; Sanchez et al. 2008; Zhang et al. 2012; Yin et al. 2018), or hybrid (Chen et al. 2018). In all instances, the surrogates in the ensembles were designed to have low bias (accurate) and to be diverse by adopting different learning algorithms (*heterogeneous ensembles*). For example, considering static integration, Acar and Rais-Rohani (2009) used ensembles of five different surrogates, namely, polynomial regression (PRS), radial basis functions (RBF), kriging (KR), Gaussian process (GP), and support vector regression (SVR). Audet et al. (2018), on the other hand, introduced an order-based error tailored to surrogate-based search, considering the PRS, RBF, and kernel smoothing (KS) models. In contrast, Yin et al. (2018) and Chen et al. (2018) reported ensembles with dynamic integration combining the typical metamodels, namely, PRS, RBF, and KR. Heterogeneous ensembles, tough, and lack of control of diversity with previous studies of ensembles in engineering are not explicitly testing the level of correlation among the surrogates in the *ensemble generation* phase. To avoid this

shortcoming, a variety of approaches have been proposed, most notably, the so-called overproduce-and-choose with pruning process (Rooney et al. 2004; Mendes-Moreira et al. 2012) and as discussed by Bishop (1995) by accounting for the correlation among the members of the ensemble (Viana et al. 2014).

Alternatively, in homogeneous ensembles (where all the surrogates share the same learning algorithm), diversity can be more systematically controlled by randomly sampling the training data (implicit methods) or deterministically manipulating the ensemble generation process (explicit methods). Randomly sampling the data so that the surrogate models are trained with different sets of data can be very effective by, for example, sampling with replacement, i.e., bagging (Breiman 1996), and the subsampling of the training set or training on random samples of input features, i.e., random subspace (Ho 1998). On the other hand, diversity can be promoted explicitly by sequentially (Rosen 1996) or simultaneously (Liu and Yao 1999) generating models which deliberately penalize the correlation between members of the ensemble in the loss function. In addition, an ensemble can be

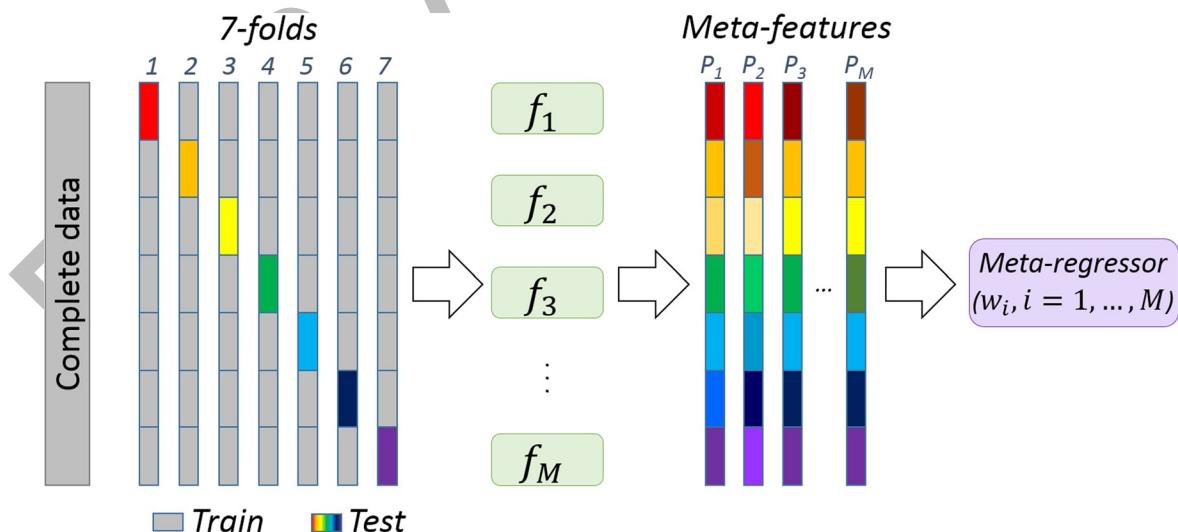
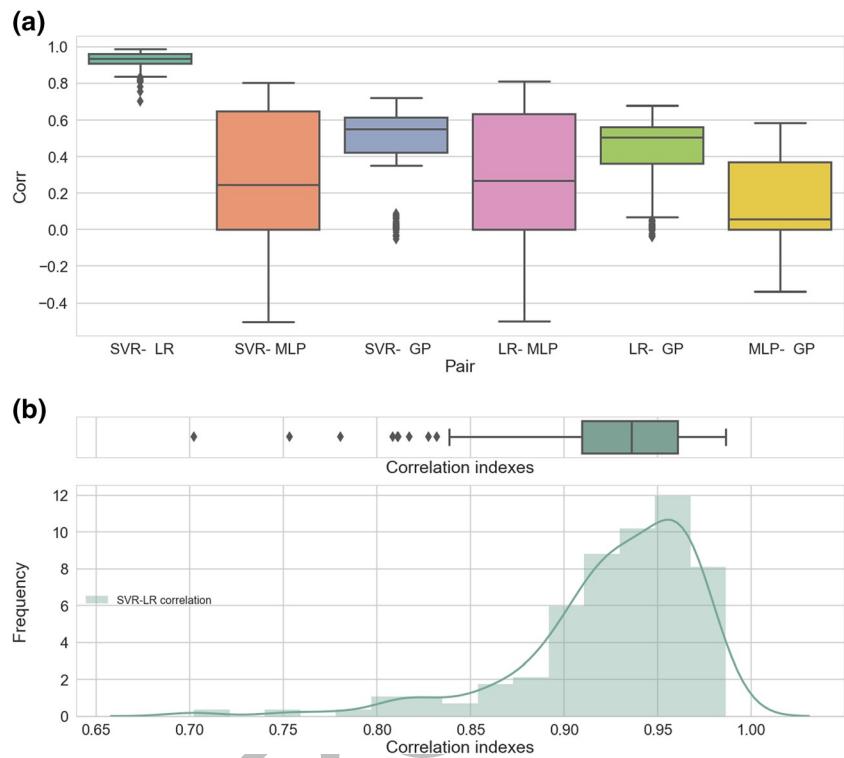


Fig. 2 Heterogeneous ensemble generation process illustrated with 7-fold cross-validation and M surrogates to construct a metalevel data (out-of-fold) set of predictions. Once the metalevel data set is available a linear

regressor can be fitted; the regression coefficients represent the relative contribution of the individual surrogates to the ensemble prediction

Fig. 3 Correlations among models of the ensemble
(a) and correlations of SVR-LinearRegression, the two most correlated models
(b)—Rosenbrock 6D test function. In this context, the labels SVR, LR, MLP, and GP, refer to support vector, linear, multilayer perceptron, and Gaussian process regressors, respectively



constructed sequentially as an additive model where at each stage a simple surrogate (boost) is fit to the negative gradient of the loss function at the training data such that a particular cost function is minimized; for example, for a quadratic loss function, the cited gradient can be shown to be the residual of the latest ensemble. Note that the contribution (weight) of each base learner to the prediction is always dynamically established, that is, it varies with location. This latest approach is known as gradient boosting (Schapire 1990; Freund 1995), a form of functional gradient descent (Breiman 1999; Friedman 2001).

Furthermore, when using ensembles and, in general, in surrogate modeling, it is often of interest to address the so-called feature selection problem. The feature selection problem makes reference to the removal of irrelevant or redundant (nonpredictive information) inputs in the training data for faster and more effective learning. More precisely, features with higher ranking or importance scores may be chosen for further investigation or for building a more parsimonious model. A frequently adopted approach for feature selection is the Sobol method (Sobol 1993), a variance-based technique used with a trained model (as a separate task) that determines the contribution of each feature and their interactions to the overall variance of the target variable. Notably, a particular form of surrogates known as decision trees can provide importance measures to preliminary rank input variable features as a by-product of the modeling process (Karagiannopoulos et al. 2007; Kazemitarab et al. 2017).

This paper presents (i) the bias/variance/covariance formulation to raise awareness of the need to address potential correlation among surrogates in heterogeneous ensembles, (ii) an ensemble approach based on gradient boosting and decision trees with diversity promoting measures for solving engineering modeling problems, and (iii) the relative performance of the proposed approach with respect to heterogeneous ensembles on selected case studies. The remainder of this paper is structured as follows: problem statement (Section 2), bias/variance/covariance formulation (Section 3), heterogeneous ensembles and correlation among surrogates (Section 4), decision trees for regression (Section 5), a gradient boosting approach for the ensemble of surrogates (Section 6), case studies (Section 7), results and discussion (Section 8), and conclusions and recommendations (Section 9).

2 Problem statement

The problem of interest is one of typical regression, that is, given a training set D consisting of n instances, such as $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, where x_i is an instance of a vector of features and y_i is a target; the task is to learn an approximate function $\hat{f} : X \rightarrow \mathbb{R}$ of the true function f from D . The approximate function \hat{f} can be a single surrogate f_i or the result of a linearly weighted sum of M individual surrogates F . The

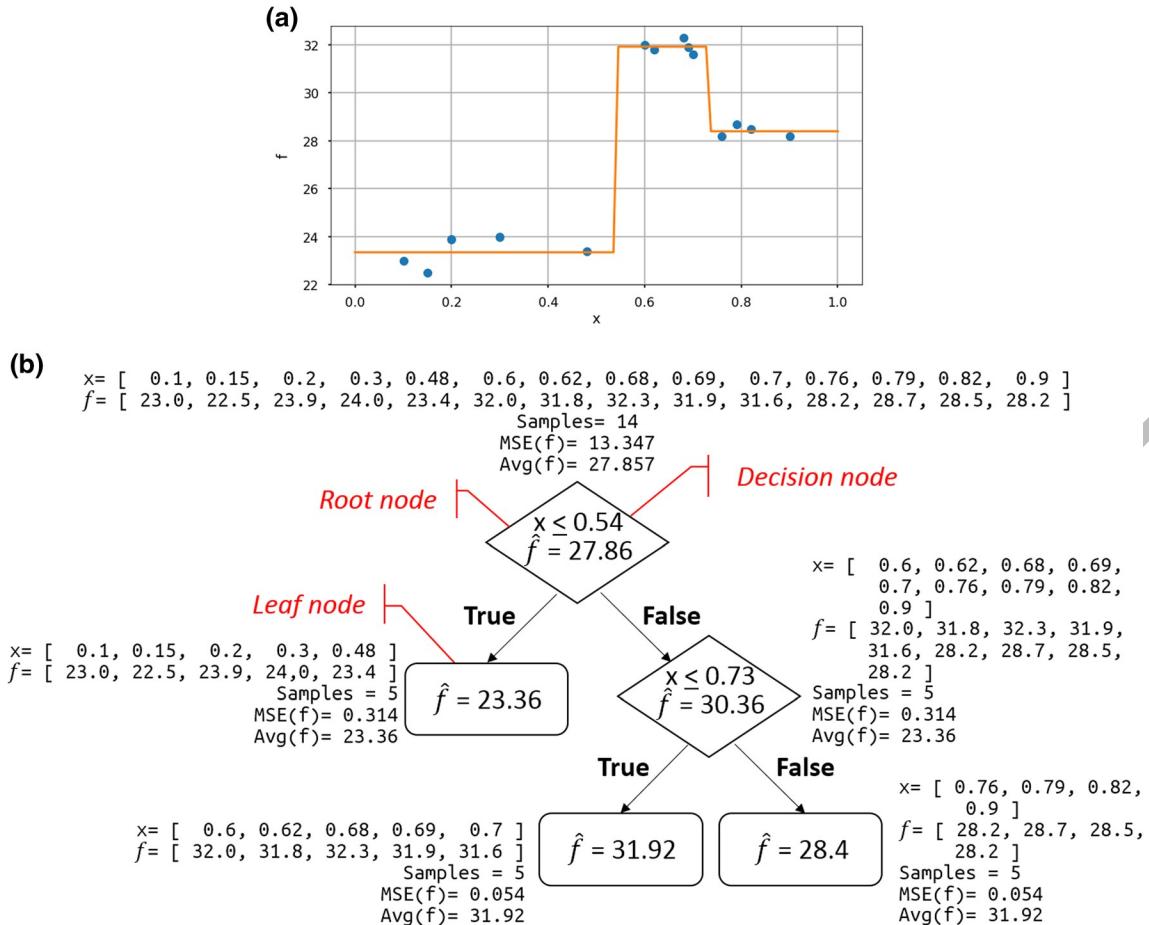


Fig. 4 A decision tree estimation from training sample (a) and the corresponding decision tree is also shown (b). Note root, decision, and leaf nodes (predictions)

function \hat{f} is estimated by optimizing a loss function typically mean square error considering out-of-fold predictions.

3 Bias, variance, and covariance decomposition

For an ensemble F (with constant weights) it can be shown (Ueda and Nakano 1996) that the expected value of the MSE-mean squared error (generalization error) can be expressed as Eq. (1):

$$E\{(F-y)^2\} = \overline{\text{bias}}^2 + \frac{1}{M} \overline{\text{var}} + \left(1 - \frac{1}{M}\right) \overline{\text{covar}} \quad (1)$$

where

$$\overline{\text{bias}} = \frac{1}{M} \sum_i (E\{f_i\} - y)^2 \quad (2)$$

$$\overline{\text{var}} = \frac{1}{M} \sum_i E\{(f_i - E\{f_i\})^2\} \quad (3)$$

$$\overline{\text{covar}} = \frac{1}{M(M-1)} \sum_i \sum_{j \neq i} E\{(f_i - E\{f_i\})(f_j - E\{f_j\})\} \quad (4)$$

where, var denotes the variance of the individual surrogates in the ensemble, and covar the averaged covariance of the surrogates in the ensemble; the estimates of all the components are computed considering several random trials (training sets). Note that in this context, diversity between surrogates (differences in the predictor outputs) is measured by their pairwise covariance. This decomposition tells us that the generalization error depends not only on the bias and variance of the individual surrogates but also on the covariance between the surrogates. The sought ensemble (optimum generalization error) should account for all components, deliberately promoting diversity among surrogates while minimizing bias.

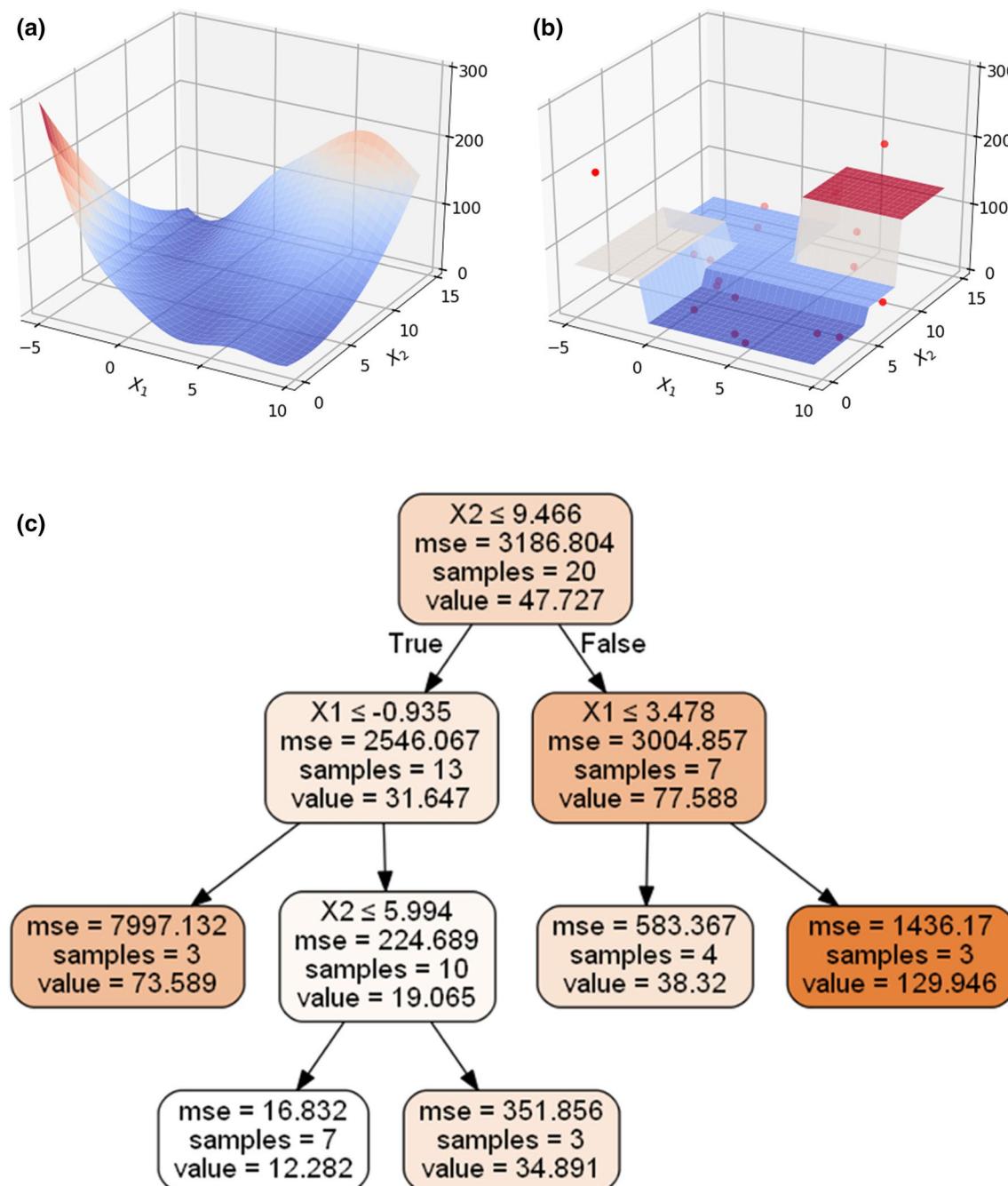


Fig. 5 Branin and Hoo 2D modeling using a decision tree. **a** Original test function. **b** Decision tree-weak learner approximation. **c** Corresponding decision tree-weak learner

Figure 1 reports the generalization error and bias, variance, and covariance decomposition for an heterogeneous ensemble of surrogates designed to approximate the well-known Rosenbrock 6D test function (Section 6), considering sixty (60) training data points and an one hundred and fifty (150) random trials. Details of the ensemble approximation are

available in the next section (Section 4). As expected, the ensemble provides a lower generalization error than the average of the individual surrogates through a reduction of the error variance. Most notably, the often ignored covariance component is shown as a prominent contributor to the generalization error, only second to the bias error.

Table 8 Hyperparameter grid for the heterogeneous ensemble—oil production NPV case study

Model	Hyperparameter	Values
Linear regression	—	—
SVR-radial basis function	C*	$135 \times 10^6, 270 \times 10^6, 540 \times 10^6$
	Gamma	0.10, 0.15, 0.20
Multilayer perceptron	Hidden layer size**	4, 6, 8
	Activation function	ReLU, logistic
Gaussian process	—	—

*The C and ϵ values under consideration are 0.5, 1.0 and 2.0 times the reference values calculated as proposed by Cherkassky and Ma (2004)

**Settled by limiting the number of parameters of the model (weights) to be a fraction (40%, 60%, 80%) of the total number of data points available to the learning process

diversity, subsampling makes reference to fitting the base learners with a fraction (subsample value) of the training data, and the random subspace value specifies the number of features to consider when looking for the best split. The hyperparameters are selected from a grid of prespecified values (grid search) using k -fold cross-validation.

II. Ensemble generation. Figure 6 shows the basic steps of using the gradient boosting algorithm with a quadratic loss function. Note it builds the ensemble by sequentially fitting a base learner (each function increment is called a boost), e.g., a *weak* decision tree, to the residuals of the latest ensemble ($y_i - F_{m-1}(x_i)$), that is, the negative gradient of the loss function at the training data (explicitly promoting diversity). *Weak* decision trees, also called *stumps*, refer to the models where a rough approximation is sought and often translates to limiting the trees depth to a very low number (often 2 to 3).

While in this study the base learners are decision trees and the loss function is quadratic, gradient boosting is not limited to a particular surrogate or loss function. More generally, at each iteration, the base learners h_m , are trained on the so-called pseudo-residuals, which represent the negative gradient Eq. (6) of the loss function at training points.

$$\nabla L = \frac{\partial L(y_i, F_{m-1}(x_i))}{\partial F_{m-1}(x_i)} \quad (6)$$

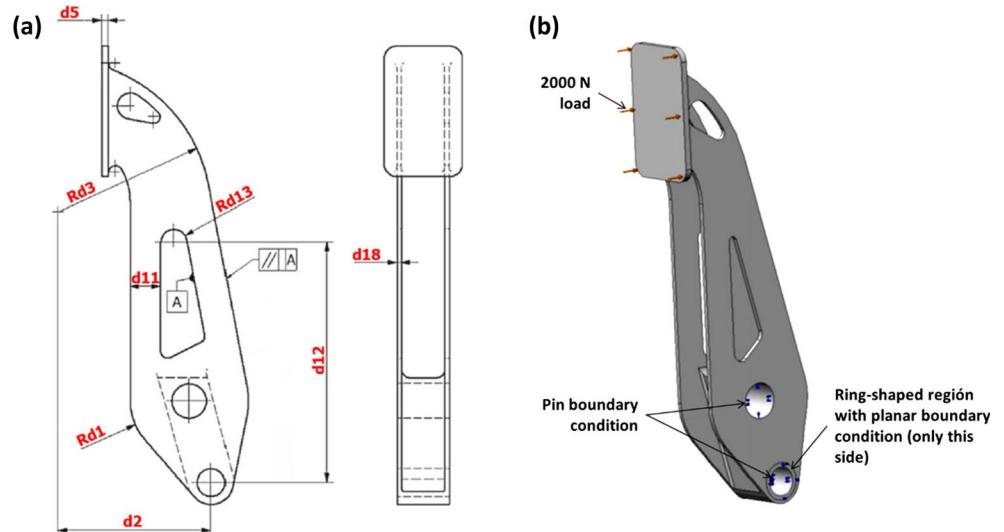
That is, the loss reduction for each training point if the predicted value— $F_{m-1}(x_i)$ —was to become one unit closer to the target value. While the base learner loses accuracy at the training points, it provides an approximation of the cited gradient throughout the input space— $h_m(x)$ —which is used to implement a greedy stagewise strategy for calculating the latest ensemble— $F_m(x)$ —i.e., $F_m(x) = F_{m-1}(x) + \alpha_m \cdot h_m(x)$. Finally, at each leaf, the optimal parameter α is sought so that $L(y_i, F_m(x_i))$ is minimized.

While the number of base learners are a parameter to be set, the gradient boosting ensembles have been shown to be resistant to overfitting so it can be specified quite robustly without jeopardizing predictive performance. An efficient and scalable implementation of the gradient boosting framework known as extreme gradient boosting, aka XGBoost, is also available (Chen and Guestrin 2016).

III. Ensemble selection. Considering the set of hyperparameters in (I), the best ensemble is the one with the lowest score in holdout or nested cross-validation (Cawley and Talbot 2010). When holdout (data split intro training and test sets) is not used, nested cross-validation should be implemented so that the same data is not used to both tune the ensemble hyperparameters and evaluate the ensemble performance, leading to overfitting and underestimation of the prediction errors. Nested CV effectively uses a series of train/validation/test set splits, where in the inner loop, ensembles are fitted to training data, and the best set of hyperparameters is the one associated with the ensemble with the lowest score over the validation set. In the outer loop, the best ensemble is identified as the one with the lowest (average) estimated generalization error over test sets. The particular approach to ensemble selection should be identified depending on the data size, model stability, and among other factors.

IV. Feature importance. Feature importance refers to the reduction in model accuracy when a variable is removed, and it is critical for variable selection and the building of the parsimonious and interpretable models. As previously mentioned, decision trees can naturally couple the training process with preliminary feature importance identification at no additional computational expense. In this study, the overall importance score for a particular feature is the average of the ones obtained in the base learners (stumps) in the ensemble.

Fig. 9 Geometric variables (a) and loads and boundary conditions (b)—FSAE brake pedal case study



8 Case studies

The proposed approach is evaluated using the three analytical test functions, namely, Branin and Hoo 2D (Dixon and Szegö 1978), Rosenbrock 6D (Rosenbrock 1960), and Dixon and Price 12D (Dixon and Price 1989) and two industrial case studies in the areas of engineering modeling, i.e., the von Mises stress and buckling load factor in a Formula SAE¹ brake pedal (structure analysis) as discussed in Romero and Queipo (2017) and the net present value of the 20-year oil production of a mature oil reservoir considering alternative new well locations (reservoir engineering).

8.1 General considerations

Table 1 lists the ensembles and individual surrogate models (base learners) under consideration and the corresponding computer implementations (libraries). The ensembles are identified as gradient boosting with diversity measures (GBwDM-proposed approach) and heterogeneous ensemble (HEA); the individual members of the heterogeneous ensemble are also indicated, i.e., linear regression, SVR-radial basis function, neural networks, and Gaussian process. All the results were obtained using the cited libraries within Python programs in Google's Colaboratory environment—a free Jupyter notebook environment than runs entirely in the Cloud. Jupyter is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations, and narrative text (www.jupyter.org).

¹ Formula SAE is a student design competition organized by SAE International, previously known as the Society of Automotive Engineers, SAE.

All the models were trained/tested using the holdout cross-validation approach; the general hyperparameters (Table 2) of all the models were identified using 5-fold cross-validation within the training set, and their values were the result of a preliminary study and adopting common practices. For all case studies, the size of the training sets was specified as 10 times the number of dimensions, and seeking sample-independent results, one hundred and fifty (150) random trials of the training/test sets were conducted.

The results of the proposed approach and those of heterogeneous ensembles are compared considering median of the mean square error in the test data set; the number of times each particular type of ensemble exhibited a lower mean square error (winner) is also reported.

8.2 Analytical test functions

The selection of these functions is aimed at providing a variety of function behavior, number of features, and ranges; Branin and Hoo (2D), Rosenbrock (6D), and Dixon and Price (12D) all frequently used as modeling and optimization case studies (Table 3). Figure 7 illustrates the two-feature behavior of the analytical test functions under consideration.

8.2.1 Training and test data sets

For all the analytical case studies, data is generated using a Latin hypercube design for both training and test sets, with sizes $10 \times$ and $100 \times$ the corresponding input domain dimension, respectively. For example, for the Rosenbrock 6D case study, the training and test sets are of sizes 60 and 600, respectively. On a separate note, Table 4 shows the hyperparameter grid under consideration for the heterogeneous ensemble.

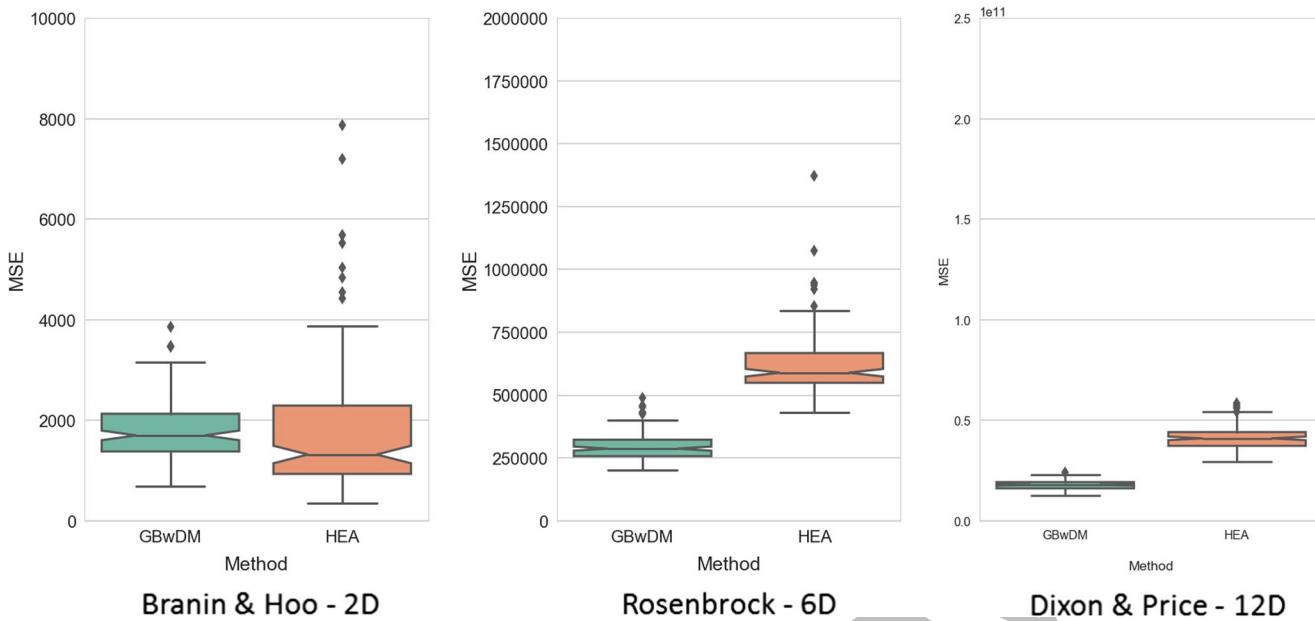


Fig. 10 Test set MSE distributions for the analytical case studies, i.e., Branin and Hoo—2D, Rosenbrock—6D, and Dixon and Price—12D

research, design, manufacturing, testing, developing, marketing, management, and finances. In particular, each component in a single-seater FSAE vehicle must be safe and lightweight in order to decrease the vehicle mass and take the most advantage of the engine power. The brake pedal design of FSAE vehicles represents an easy-to-understand structural problem among engineering students and practitioners around the world yet relevant to aerospace and automotive industrial environments.

With reference to Fig. 9a, the problem of interest is to estimate the von Mises stress (ε_S) and buckling load factor (ε_B) in a FSAE brake pedal considering as features eight (8) geometric variables (in circles) and the Young modulus E of the material. The loads and boundary conditions of the brake pedal are also shown in Fig. 9b with the range of the features specified in Table 9. Please see Romero and Queipo (2017) for a detailed discussion.

8.5 Training and test data sets

The data consists of 1320 samples randomly divided in training and test data sets of sizes 90 and 1230, respectively.

Table 10 shows the hyperparameter grid for the heterogeneous ensemble.

9 Results and discussion

The proposed approach (GBwDM) exhibited the following:

- The lowest median (statistically significant) of test set mean square errors and most winners for analytical (except the 2D) and some industrial case studies, specifically, oil production NPV optimization and FSAE – von Mises stress
- Strong resistance to overfitting for all case studies

Table 11 Median of test set MSE for the analytical case studies, i.e., Branin and Hoo—2D, Rosenbrock—6D, and Dixon and Price—12D

MSE (test set)	Branin and Hoo 2D	Rosenbrock 6D	Dixon and Price 12D
GBwDM	1538.59	283, 234.06	19.77×10^9
HEA	681.21	366, 008.06	20.96×10^9
Pearson's chi-squared test (p value)	<< 0.05	<< 0.05	<< 0.05

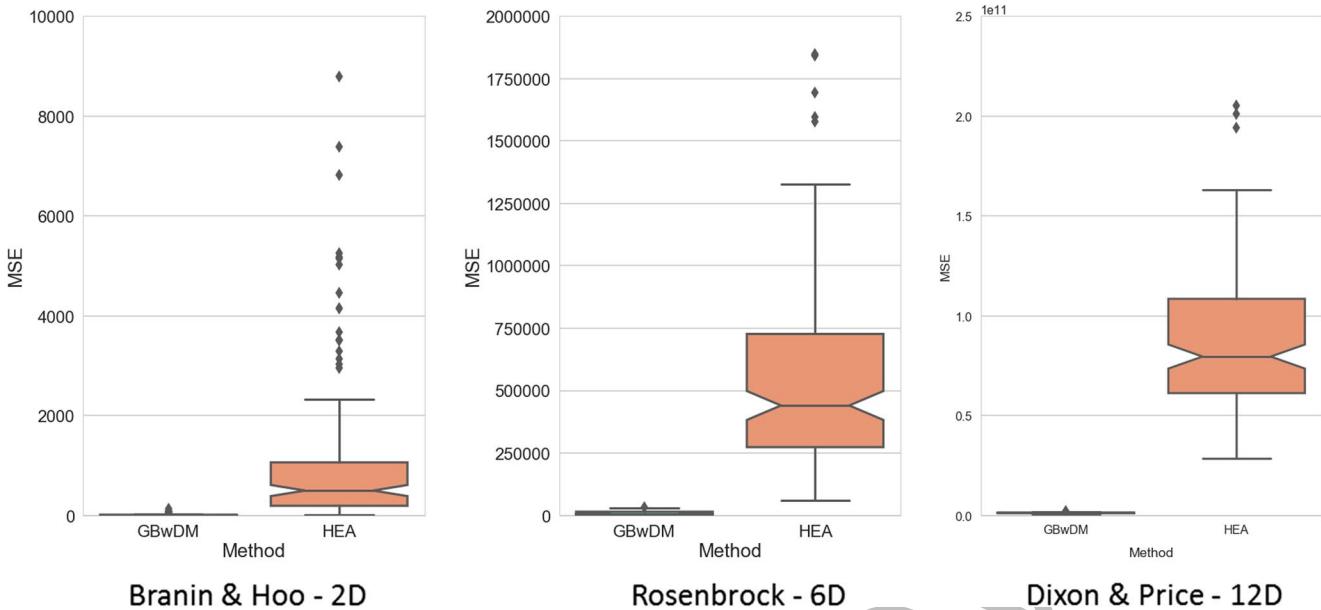


Fig. 11 Training MSE distributions for the analytical case studies, i.e., Branin and Hoo—2D, Rosenbrock—6D, and Dixon and Price—12D

- Considerable speedup with respect to heterogeneous ensemble learning
- Diversity measures that proved to be effective for improving the prediction performance
- A preliminary assessment of feature importance that is statistically aligned with those reported by the frequently used Sobol method

The lowest median of test set mean square errors for all analytical (except the 2D) case studies. Based on one hundred and fifty (150) random trials, GBwDM outperformed the heterogeneous ensemble approach (HEA) with 56% (6D) and 71% (12D) lower median MSE (Fig. 10; Table 11); the error distribution associated with the training data is shown in Fig. 11. Pearson's chi-squared test on the difference of the

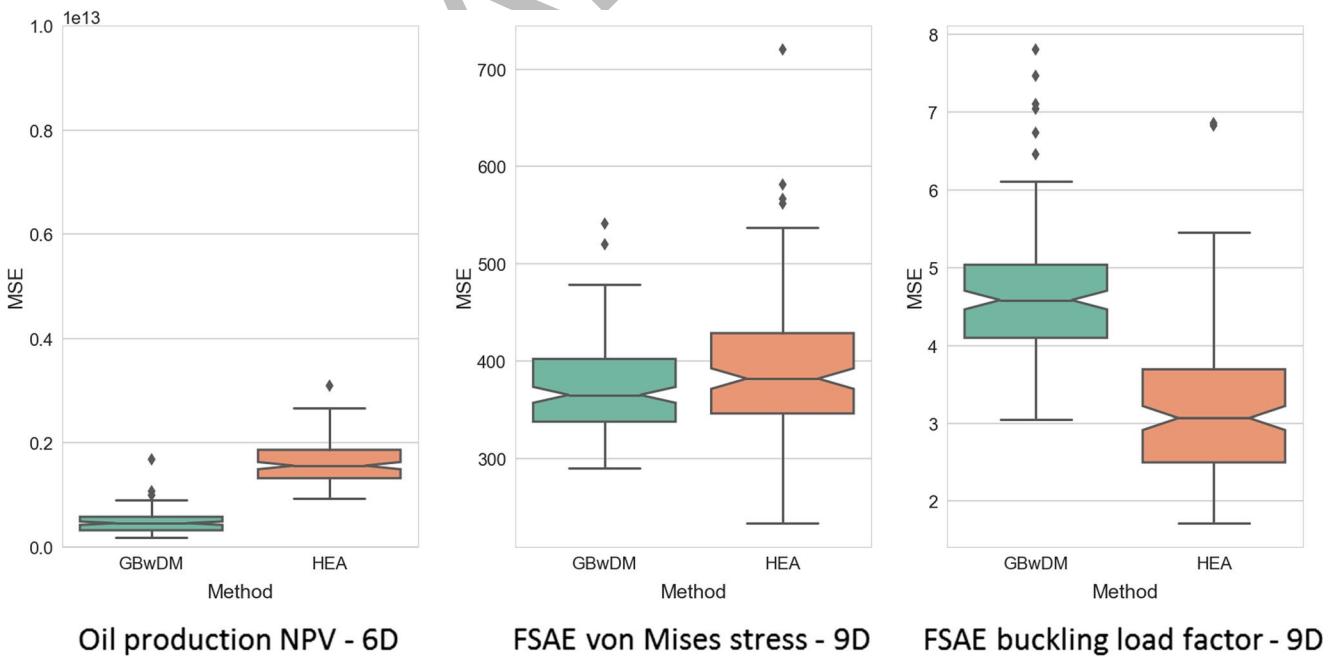
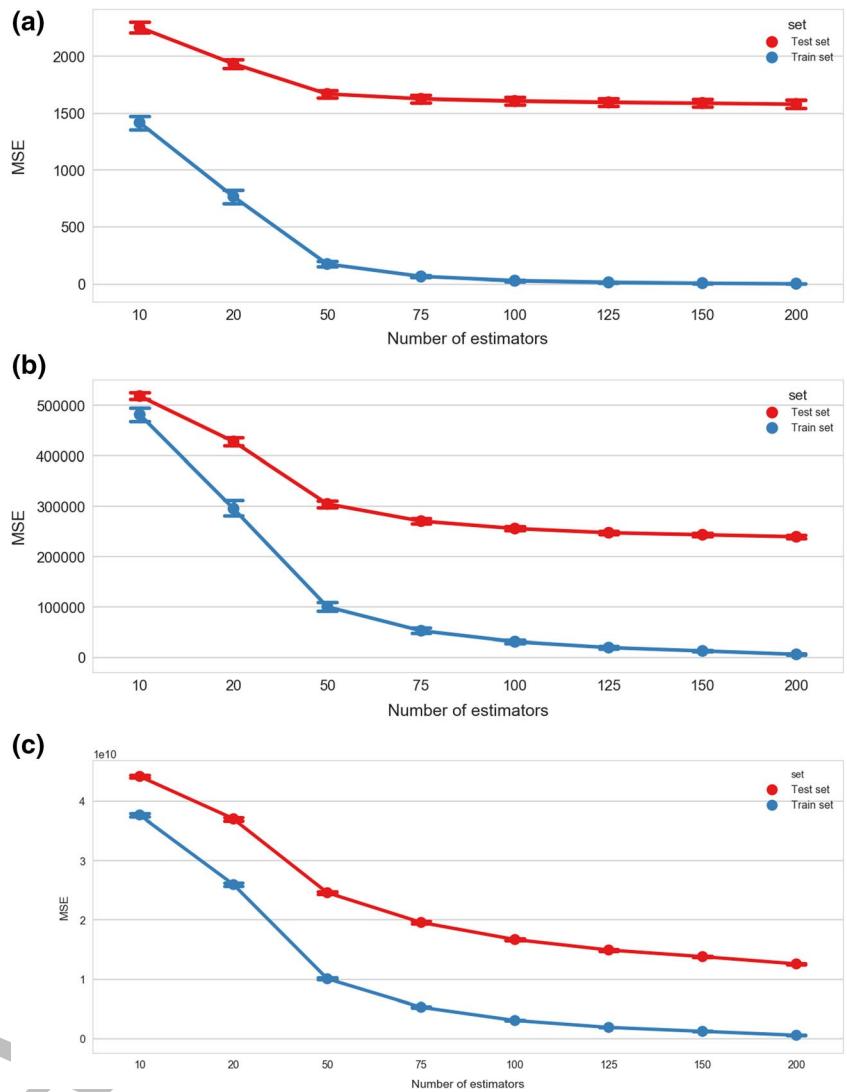


Fig. 12 Test set MSE distribution for the industrial case studies, i.e., oil production NPV, FSAE-von Mises stress, and FSAE-buckling load factor

Fig. 15 Median of training and test set MSE for the analytical case studies, i.e., Branin and Hoo—2D (a), Rosenbrock—6D (b), and Dixon and Price—12D (c)



random trial of the training set in the FSAE-von Mises stress and FSAE-buckling load factor studies, respectively. Note that for the former target variable, the three (3) most important features (with reference to Fig. 17a) are $d18$ (69%), $d1$ (12.1%), and $d12$ (5.2%), while for the latter (Fig. 18a), those are $d18$ (95.2%), the module of elasticity E (2.1%), and $d12$ (1.4%); in parenthesis, the reduction percentage of the total MSE. Considering the learning algorithm for decision trees greedily select the splitting variables for maximum reduction of a variability measure (e.g., MSE), one would assume that the feature ranking should be related to that of the Sobol method. That indeed seems to be the case. For one hundred and fifty (150) random trials of the training set, Fig. 17 b and Fig. 18 b show the empirical distribution of a measure of the rank correlation-similarity (Kendall's tau statistic) of the ordering of the most important features obtained from GBwDM and the Sobol method for both target variables of interest. The median of the corresponding p values confirms the statistical

dependence of the feature rankings from GBwDM and the Sobol method (Fig. 17c; Fig. 18c). Specifically, for the target variables of interest, the medians of the p values are in both cases approximately 1×10^{-3} , significantly lower than those of the threshold of 0.05 for rejecting the null hypothesis (the two ranking features are independent).

10 Conclusions and recommendations

This work presented a homogeneous ensemble approach based on gradient boosting and decision trees with diversity measures (subsampling and random subspace) for solving engineering modeling problems. The proposed approach is designed to improve prediction estimates by promoting diversity among the members of the ensemble (reducing the covariance component of the generalization error), be robust to overfitting, and, as a

Fig. 16 Median of training and test set MSE for the industrial case studies, i.e., oil production NPV—6D (a), FSAE-von Mises stress—9D (b), and FSAE-buckling load factor—9D (c)

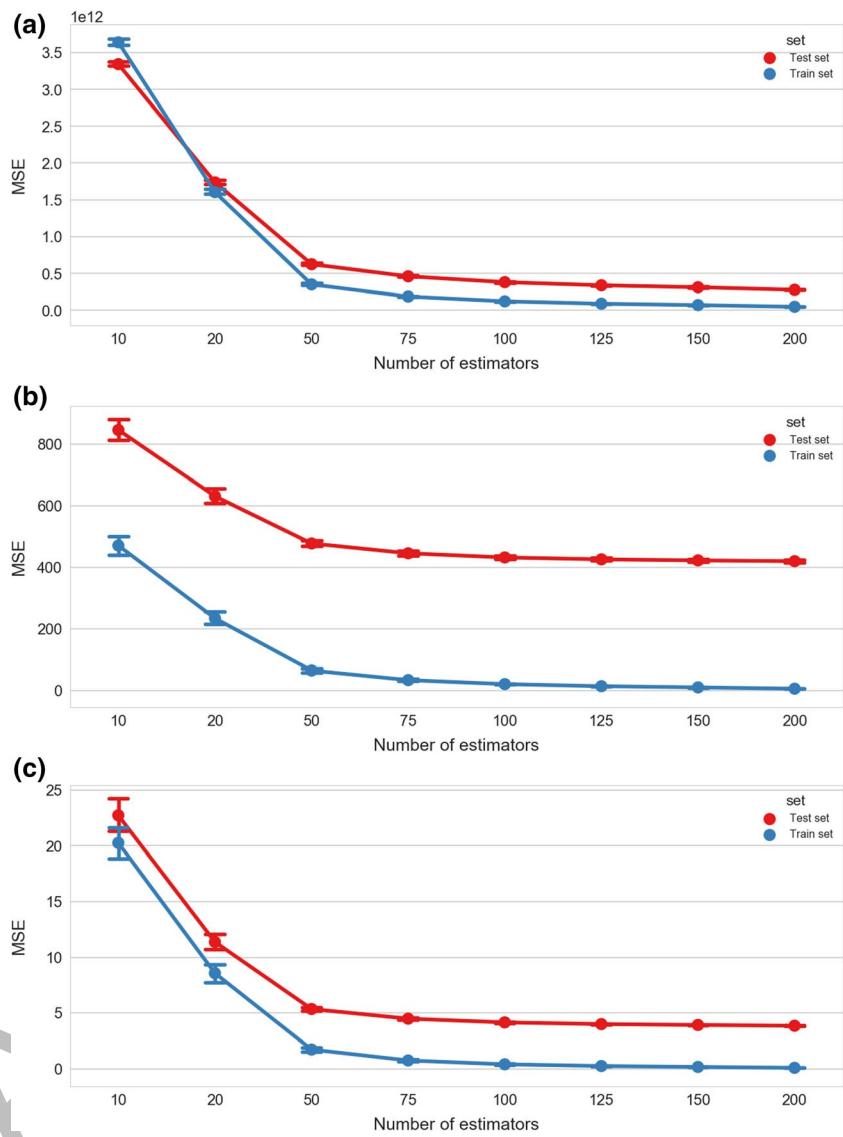
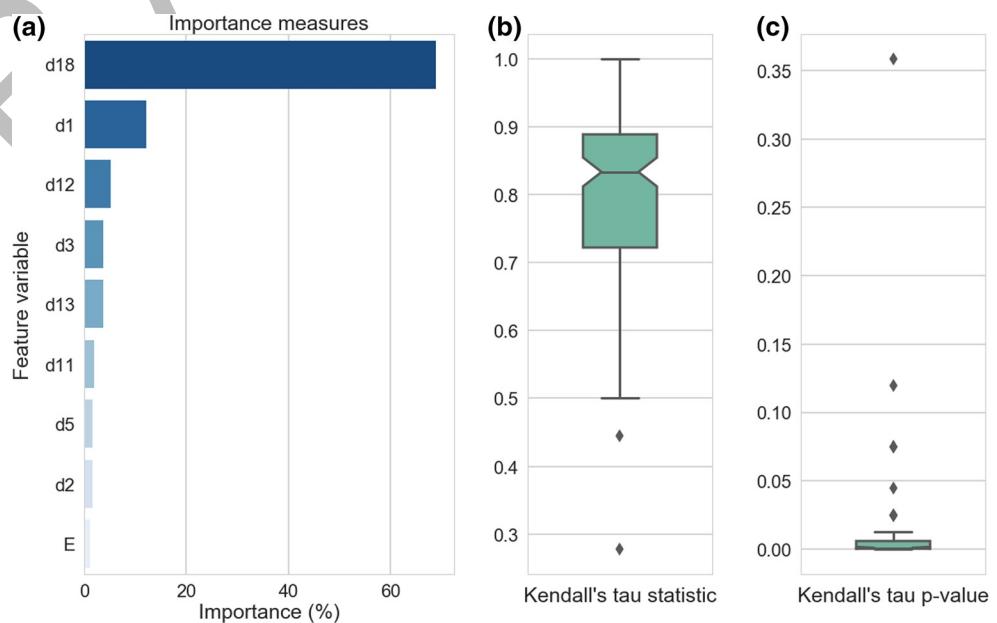


Fig. 17 GBwDM feature importance for a selected random trial (a), a rank correlation-similarity (Kendall's tau) statistic distribution (b), and the corresponding p value (c), for one hundred and fifty (150) random trials when compared with Sobol's feature importance ranking—FSAE-von Mises stress case study



- Sobol I (1993) Sensitivity analysis for non-linear mathematical models. *Mathematical Modeling & Computational Experiment* 1:407–414
- Syed A (2012) Technology focus: mature fields and well revitalization. *J Pet Technol* 64:74–74. <https://doi.org/10.2118/0112-0074-JPT>
- Tenne Y (2013) An optimization algorithm employing multiple metamodels and optimizers. *Int J Autom Comput* 10(3):227–241. <https://doi.org/10.1007/s11633-013-0716-y>
- Ueda N and Nakano R (1996). Generalization error of ensemble estimators. In Proceedings of International Conference on Neural Networks(ICNN'96), 90–95. <https://doi.org/10.1109/ICNN.1996.548872>
- Viana FAC, Simpson TW, Balabanov V, Toropov V (2014) Metamodelling in multidisciplinary design optimization: how far have we really come? *AIAA J* 52(4):670–690. <https://doi.org/10.2514/1.J052375>
- Yin H, Fang H, Wen G, Gutowski M, Xiao Y (2018) On the ensemble of metamodels with multiple regional optimized weight factors. *Struct Multidisc Optim* 58:245–263. <https://doi.org/10.1007/s00158-017-1891-1>
- Zerpa L, Queipo NV, Pintos S, Salager JL (2005) An optimization methodology of alkaline–surfactant–polymer flooding processes using field scale numerical simulation and multiple surrogates. *J Pet Sci Eng* 47(3–4):197–208. <https://doi.org/10.1016/j.petrol.2005.03.002>
- Zhang J, Chowdhury S, Messac A (2012) An adaptive hybrid surrogate model. *Struct Multi Optim* 46(2):223–238. <https://doi.org/10.1007/s00158-012-0764-x>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.