

# A Novel Sensitivity-based Method for Feature Selection

Dayakar L. Naik<sup>1</sup> and Ravi kiran<sup>2</sup>

<sup>1</sup> Research Associate, Department of Civil & Environmental Engineering, North Dakota State University, ND 58105, email: [dayakarnaik.lavadiya@ndsu.edu](mailto:dayakarnaik.lavadiya@ndsu.edu).

<sup>2</sup> Assistant Professor (Corresponding Author), Department of Civil & Environmental Engineering, North Dakota State University, ND 58105, email: [ravi.kiran@ndsu.edu](mailto:ravi.kiran@ndsu.edu).

## Abstract

Sensitivity analysis is a popular feature selection approach employed to identify the important features in a dataset. In sensitivity analysis, each input feature is perturbed one-at-a-time and the response of the machine learning model is examined to determine the feature's rank. Note that the existing perturbation techniques may lead to inaccurate feature ranking due to their sensitivity to perturbation parameters. This study proposes a novel approach that involves the perturbation of input features using a complex-step. The implementation of complex-step perturbation in the framework of deep neural networks as a feature selection method is provided in this paper, and its efficacy in determining important features for real-world datasets is demonstrated. Furthermore, the filter-based feature selection methods are employed, and the results obtained from the proposed method are compared. While the results obtained for the classification task indicated that the proposed method outperformed other feature ranking methods, in the case of the regression task, it was found to perform more or less similar to that of other feature ranking methods.

**Keywords:** Complex step derivative approximation (CSDA); Feature Ranking; Regression; Classification; Feature relevance; and neural networks.

## 1. Introduction

Feature selection is a process of identifying a subset of features that dictate the prediction accuracy of the target variables/ class labels in a given machine learning task [1–3]. Identification of relevant features improves the machine learning (ML) models' generalized performance and facilitates a better understanding of the data in relation to the ML model [4]. For performing the task of feature selection, various methods have been proposed by researchers in the past. These methods could be broadly grouped into six categories, namely, filter methods, wrapper methods, embedded methods, hybrid methods, ensemble methods, and integrative methods [5–7]. While filter methods select features based on a performance metric regardless of the supervised learning algorithm [8–12], the wrapper methods choose feature subset by iteratively examining a certain or an ensemble of the ML algorithm's performance for selected features [13]. Examples of filter methods include Pearson correlation coefficient, information gain, gain ratio, Chi-square, Fisher score, ReliefF, etc., and examples of wrapper method include sequential feature selection, genetic algorithms, etc. On the other hand, in embedded methods, the feature selection algorithm is integrated into the learning algorithm [5,9,13]. Examples of the embedded method include decision tree, random forest, support vector machine recursive feature elimination (SVM-RFE). When compared to filter-based approaches, the embedded approach yields higher accuracy because of its interaction with a specific classification model. A comprehensive review of these three methods' description and comparison is discussed by various researchers in the literature [4,5,14–19].

In hybrid methods, multiple conjunct primary feature selection methods are applied consecutively [6]. For instance, Liu et al. [20] proposed a hybrid feature selection method in which mutual information was first applied to identify the relevant features from the feature set, and then the wrapper method was applied subsequently to choose the subset of best features from the relevant features. Ensemble feature selection methods use an aggregate of feature subsets of diverse base classifiers [6]. For instance, Hoque et al. [21] proposed an Ensemble Feature Selection – Feature Selection (EMI-FS) in which information gain, gain ratio, ReliefF, symmetric uncertainty, and Chi-square were employed as base filter methods to obtain the relevant subset of features which were subsequently combined to extract the optimal subset. In the integrative feature selection method, the external knowledge of feature selection is integrated [6]. For example, Cindy et al. [7], proposed an integrative gene selection approach in which gene rankings are determined by considering both the statistical significance of a gene in the dataset and the biological background information acquired through research. In this paper, we restrict our scope to the embedded feature selection methods that incorporate feed-forward neural networks/multi-layer perceptron as the learning models.

Multi-layer Perceptron (MLP) is a basic type of neural network that learns a function  $g: \mathbb{R}^q \rightarrow \mathbb{R}^m$  by training on a dataset, where  $q$  is the number of inputs and  $m$  is the number of outputs. MLP's were employed for performing feature selection by various researchers in the past. For instance, Setiono and Liu [22] developed a neural network feature selector method based on backward elimination wherein weights of low magnitude in the network were converged to zero by adding a penalty term to the error function. Sindhwani et al. [23] presented a maximum output information algorithm for feature selection. Liefeng Bo [24] proposed MLP Embedded Feature Selection

(MLP-EFS), in which each feature is multiplied by the corresponding scaling factor. By applying truncated Laplace prior to the scaling factors, feature selection is integrated into MLP-EFS.

Notwithstanding to methods mentioned above, sensitivity analysis of MLP and support vector machines (SVM) was also carried out to perform feature selection. For instance, Ruck et al. [25] developed a technique that analyzes the weights in MLP to determine essential features. Gasca et al. [26] proposed a saliency measure that estimates the input features' relative contribution to the output neurons. Utans et al. [20] proposed a 'sensitivity-based-pruning (SBP)' to remove irrelevant input features from a nonlinear regression model. Acir et. al. [29] implemented the perturbation method in the framework of SVM to perform feature selection for classification of Electrocardiogram (ECG) beats. Sensitivity analysis examines the change in the target output when one of the input features is perturbed, i.e., first-order derivatives of the target variable with respect to the input feature are evaluated. Herein we refer the first-order derivative term as the feature sensitivity metric. The higher the magnitude of change in feature sensitivity metric, the higher is the importance of input feature. At this juncture, it is important to note that sensitivity analysis methods involve computation of the feature sensitivity metric or first-order derivative for identifying important features. In general backpropagation algorithm (for MLP), is employed or finite difference schemes [30–33] is used for computing feature sensitivity metric. Employing numerical differentiation techniques such as finite difference approximation (FDA) (see Eq. 1) and central finite difference approximation (CFDA) (see Eq. 2) results in inaccurate computation of derivatives [34,35] because of inappropriate choice of step size. For instance, Juana et.al. [36] introduced the iterative perturbation method for auto-tuning the step size for SVM. Such errors arising due to the choice of smaller step sizes are referred to as subtractive cancellation errors.

Finite difference approximation (FDA)

$$g'(x_1, x_2, \dots, x_k, \dots, x_q) \approx \frac{f(x_1, x_2, \dots, x_k + h, \dots, x_q) - f(x_1, x_2, \dots, x_k, \dots, x_q)}{h} \quad (1)$$

Central finite difference (CFDA)

$$g'(x_1, x_2, \dots, x_k, \dots, x_q) \approx \frac{f(x_1, x_2, \dots, x_k + h, \dots, x_q) - f(x_1, x_2, \dots, x_k - h, \dots, x_q)}{2h} \quad (2)$$

where  $\mathbf{x} = (x_1, x_2, \dots, x_k, \dots, x_q)' \in \mathbb{R}^{q \times 1}$  are the inputs,  $q$  is the number of inputs,  $g(\cdot)$  is the function mapping the inputs to the output variable and,  $g'(\cdot)$  is the first partial derivative approximation of  $f(\cdot)$  with respect to the input  $x_k$ . The feature  $x_k$  is perturbed in both the cases to get the first derivative as seen in Eq. (1) and (2).

In this paper, a novel Complex-step sensitivity analysis-based feature selection method referred to as CS-FS is proposed, which incorporates a complex-step perturbation of the input feature to compute the feature sensitivity metric and identify the important features. It evaluates the analytical quality first-order derivatives without the need for extra computations in neural networks or SVM ML models. A brief overview of the complex step perturbation approach is provided in Section 2, and its implementation in the framework of FFNN to perform feature selection is described in Section 3. The details of the dataset are provided in Section 4 and the

efficacy of the proposed method is then demonstrated on real-world datasets in Section 5, and the summary and future work are provided in Section 6.

## 2. Overview of Complex-Step Perturbation Approach (CSPA)

CSPA, originally referred to as complex-step derivative approximation (CSDA), was proposed by Lyness and Moler [37] to evaluate the first-order derivative of analytic functions. A simplified version of mathematical derivation for computing the first-order derivative of a scalar function using complex-step perturbation was then provided by Squire and Trapp [38] which is as follows.

Consider a holomorphic function  $f(\cdot)$  which is infinitely differentiable. The Taylor series expansion of the function  $f(\cdot)$  evaluated at the complex perturbed point  $x_0 + ih$  is expressed as

$$f(x_0 + ih) = f(x_0) + ihf'(x_0) - \frac{h^2}{2!}f''(x_0) - \frac{ih^3}{3!}f'''(x_0) + \dots \quad (3)$$

where,  $h$  is the step size and  $i^2 = -1$ .

By taking the imaginary component of  $f(x_0 + ih)$ , and truncating the higher-order terms in the Taylor series, the first-order derivative can be expressed as

$$f'(x_0) = \frac{\text{Imag}(f(x_0 + ih))}{h} + \mathcal{O}(h^2) \quad (4)$$

where,  $\text{Imag}(\cdot)$  denotes the imaginary component and  $\mathcal{O}(h^2)$  is the second-order truncation error. It is evident from Eq. 4 that the first-order derivative evaluated using the CSPA technique is not prone to subtractive cancellation errors (see Eq.1 and Eq.2) due to the absence of subtractive operations. Furthermore, a choice of the small magnitude of  $h$  could possibly eliminate the truncation error  $\mathcal{O}(h^2)$  too. A simple example illustrating the accuracy of CSPA over finite difference schemes can be found elsewhere [39,40]. Some examples of the fields where CSPA is currently gaining a lot of attention for performing sensitivity analysis includes aerospace [41–44], computational mechanics [39,40,45], estimation theory (e.g., second-order Kalman filter) [46].

## 3. Complex-step Feature Selection Method

In the proposed method, we implement a complex-step perturbation in the framework of feed-forward neural networks to illustrate the task of feature selection. Note that this could be extended to other ML models such as SVM whose decision function is holomorphic. Higher the change in the magnitude of the output variable  $y \in \mathbb{R}$  of the FFNN with respect to the input feature  $x_k \in \mathbb{R}$ , higher is the importance of the feature  $x_k$ . For a multivariate function, the extended form of CSPA can be expressed as

$$g'(x_1, x_2, \dots, x_k, \dots, x_q) = \frac{\text{Imag}(g(x_1, x_2, \dots, x_k + ih, \dots, x_q))}{h} + \mathcal{O}(h^2) \quad (5)$$

where  $\mathbf{x} = (x_1, x_2, \dots, x_k, \dots, x_q)' \in \mathbb{R}^{q \times 1}$  is a vector of input features,  $q$  is the number of input features,  $g(\cdot)$  is the function mapping the input features to the output target variable and,  $g'(\cdot)$  is the first-order derivative approximation of  $g(\cdot)$  with respect to the  $k^{\text{th}}$  input feature  $x_k$ .

### 3.1. Feature Selection for Regression Using Complex-step Sensitivity

The proposed feature selection method for the regression task involves four steps (see Figure 1). In the first step, an FFNN is configured and trained for a given dataset. Configuring the FFNN is a trial-and-error process that involves finding the appropriate number of neurons and hidden layers in a network. A neural network is said to be configured when it is capable of learning a mathematical mapping between the input features and the associated target variable such that it could be generalized to the unseen data instances. In the second step, one of the input features,  $x_k$  is chosen at a time and is perturbed with an imaginary step size of  $ih$  (where  $h \ll 10^{-8}$ ). Feedforward operation is then performed with the perturbed feature on the trained FFNN, and the results in the output layer are obtained. In the third step, the imaginary components of the output neurons' results are extracted for each perturbed feature and are divided with the step size ( $h$ ) (see Eq. 5), i.e., the first-order derivative of the target output with respect to the input feature is evaluated. Note that step 2 and step 3 are repeated for all instances in the dataset, and the average absolute magnitude of the first-order derivative of the target output with respect to the input feature is evaluated. For example, if  $y$  is the target output variable and  $x_{jk}$  is the  $k^{th}$  feature in the  $j^{th}$  observation that is complex-step perturbed ( $ih$ ), then the first order derivative of the target output with respect to the input feature averaged over all instances of datasets is expressed as (see Eq. 6)

$$\frac{\partial y}{\partial x_k} = \frac{1}{N} \sum_{j=1}^N \left| \frac{\partial y}{\partial x_{jk}} \right| \quad (6)$$

where,  $N$  denotes the number of instances in the dataset,  $k = 1 \dots q$  indicates the input feature, and  $j$  represents the observation number in the dataset. In the fourth and final step, the rank of each input feature is determined based on the magnitude of the first-order derivatives evaluated, as shown in Eq. 5. The feature with a higher magnitude of the first-order derivative is assigned a higher rank and vice versa. Note that for training the feedforward neural network, a backpropagation algorithm, in conjunction with the Levenberg-Marquardt optimization technique, is employed in this study [47].

### 3.2. Feature Selection for Classification Using Complex-Step Sensitivity

Unlike regression, a modification to step 3 is needed in the proposed method when feature selection is performed on the classification task, i.e., evaluating the first-order derivative of target output with respect to perturbed input feature. The need for modification could be attributed to two reasons: (1) discrete output in the output layer and (2) multiple first-order derivatives yielded by the feed-forward neural network output layer (SoftMax layer) (see Figure 2). Considering the fact that the inputs fed to the SoftMax activation neurons in the output layer are not discrete, the first-order derivatives of such inputs could still be evaluated. These first-order derivatives will aid in providing information about the importance of the input features. If  $\Sigma_r$  represents the net function of  $r^{th}$  neuron in the SoftMax layer, then the first-order derivative of the net function  $\Sigma_r$  with respect to the  $k^{th}$  feature  $x_k$  is expressed as (see Eq. 7)

$$\left( \frac{\partial \Sigma_r}{\partial x_k} \right) = \frac{1}{h} \text{Imag}(\Sigma_r(x_k + ih)) \quad (7)$$

where,  $r = 1 \dots m$  and  $m$  indicates the number of class labels. To quantify the change in the target output with respect to the  $k^{th}$  input feature  $x_k$ , the average of the first-order derivatives obtained for all neurons in the output layer is determined. This average magnitude is referred to as saliency ( $S_k$ ) of  $k^{th}$  input feature [25] and is expressed as (see Eq. 8)

$$S_k = \frac{1}{N} \sum_{j=1}^N \left( \sum_{r=1}^m \left| \left( \frac{\partial \Sigma_r}{\partial x_{jk}} \right) \right| \right) \quad (8)$$

where  $r$  denotes the neuron in the SoftMax output layer,  $m$  represents the number of class labels,  $\Sigma_r$  represents the net function of  $r^{th}$  neuron in the SoftMax layer. The rank of each input feature is then determined based on the magnitude of the first-order derivatives for each perturbed feature  $x_k$  determined as shown in Eq. 8.

#### 4. Numerical Experiments

In this section, numerical experiments are performed to demonstrate the effectiveness of the proposed method.

##### 4.1. Datasets

Three real-world datasets, each for regression and classification problems, are employed to demonstrate the proposed method's efficacy. The datasets are obtained from the UCI open-source data repository [48]. For regression problems, the body fat percentage dataset, abalone dataset, and wine quality dataset are chosen, and, for the classification task, a vehicle dataset, segmentation dataset, and breast cancer dataset are chosen. One of the main reasons for choosing these datasets is that they are commonly adopted in the literature of feature selection. On the other hand, the results obtained from some of the chosen datasets such as body fat percentage, wine quality, segmentation are easily interpretable and aids in ensuring the verification of the proposed method. While most of the chosen datasets have descriptive features that are continuous in nature, the proposed method can be extended to the datasets consisting of discrete input features. The descriptive features and target variables for each dataset are mentioned as follows.

##### *Regression*

Body fat percentage dataset [49]: Features – (1) Age (years), (2) Weight (kg), (3) Height (cm), (4) Neck (cm), (5) Chest (cm), (6) Abdomen (cm), (7) Hip (cm), (8) Thigh (cm), (9) Knee (cm), (10) Ankle (cm), (11) Biceps (cm), (12) Forearm (cm), (13) Wrist (cm); Target variable – percentage of body fat.

Abalone dataset [50]: Features – (1) Female, (2) Infant, (3) Male, (4) Length (gms.), (5) Diameter (gms.), (6) Height (gms.), (7) Whole weight (gms.), (8) Shucked weight (gms.), (9) Viscera weight (gms.), (10) Shell weight (gms.); Target variable – Number of rings.

Wine quality dataset [51]: Features – (1) fixed acidity, (2) volatile acidity, (3) citric acid, (4) residual sugar, (5) chlorides, (6) free sulfur dioxide, (7) total sulfur dioxide, (8) density, (9) pH, (10) sulfates, (11) alcohol; Target variable – quality score (1 to 10).

##### *Classification*

Vehicle dataset [52]: Features – (1) Compactness, (2) circularity, (3) radius circularity, (4) radius ratio, (5) axis aspect ratio, (6) maximum length aspect ratio, (7) scatter ratio, (8) elongatedness, (9) axis rectangularity, (10) maximum length rectangularity, (11) scaled variance major, (12) scaled variance minor, (13) scaled radius of gyration, (14) skewness major, (15) skewness minor, (16) kurtosis major, (17) kurtosis minor, (18) hollow ratio; Target variable – Class label 1 (van), Class label 2 (Saab), Class label 3 (bus), Class label 4 (Opel).

Segmentation dataset [48]: Features – (1) region-centroid-col (2) region-centroid-row (3) short-line-density (4) the results of a line extraction algorithm that counts how many lines of length (5) vedge-mean (6) vedge-sd (7) hedge-mean (8) hedge-sd (9) intensity-mean (10) rawred-mean (11) rawblue-mean (12) rawgreen-mean (13) exred-mean (14) exblue-mean (15) exgreen-mean (16) value-mean (17) saturatoin-mean (18) hue-mean; Target variable – Class label 1 (Window), Class label 2 (foilage), Class label 3 (brickface), Class label 4 (path), Class label 5 (cement), Class label 6 (grass), Class label 7 (sky).

Breast cancer dataset [53]: Features – (1) radius1, (2) texture1, (3) perimeter1, (4) area1, (5) smoothness1, (6) compactness1, (7) concavity1, (8) concave points1, (9) symmetry1, (10) fractal dimension1, (11) radius2, (12) texture2, (13) perimeter2, (14) area2, (15) smoothness2, (16) compactness2, (17) concavity2, (18) concave points2, (19) symmetry2, (20) fractal dimension2, (21) radius3, (22) texture3, (23) perimeter3, (24) area3, (25) smoothness3, (26) compactness3, (27) concavity3, (28) concave points3, (29) symmetry3, (30) fractal dimension3; Target variable – Class label 1 (Benign), Class label 2 (Malignant).

Other details about regression and classification datasets are provided in Table 1 and Table 2, respectively.

## 4.2. Configuring feed-forward neural networks

Feed-forward neural networks (FFNN) with three hidden layers (HL) are configured to train on the regression and classification datasets. While a configuration of 1<sup>st</sup> HL – 20 neurons, 2<sup>nd</sup> HL – 10 neurons, and 3<sup>rd</sup> HL – 5 neurons is employed to train on regression datasets, a configuration of 1<sup>st</sup> HL – 60 neurons, 2<sup>nd</sup> HL – 40 neurons, and 3<sup>rd</sup> HL – 20 neurons is employed to train on classification datasets. A Rectified Linear Unit (ReLU) nonlinear function is used as an activation function for all the configurations [54]. Note that different architectures and model parameters yield different results if a suitable configuration is not adopted. In this study, various trail configurations of increased complexity (i.e., more hidden neurons and hidden layers) were examined before choosing a suitable configuration. Herein, the suitable configuration refers to the model architecture for which further improvement in performance was not observed with an increase in complexity of architecture. For training, validating, and testing the chosen configurations, the datasets are randomly partitioned into 70:15:15 ratio, respectively. Note that in the case of the classification task, the partition ratio is maintained consistently for each class label, i.e., 70:15:15 of training, validation, and testing data from each class label is chosen. To ensure that the chosen configurations yield repeatable results, the training operation is performed 100 times with the same partition ratio but with the replacement of instances randomly selected in every iteration. The performance metric, namely mean squared error (MSE) and accuracy, are

evaluated for regression and classification datasets, respectively, for chosen configurations. The average MSE error for body fat percentage, abalone, and wine quality datasets is determined to be 20.41, 4.6, and 0.53, respectively. The average accuracy for the vehicle, segmentation, and breast cancer dataset is determined to be 75%, 80% and, 90%, respectively. The addition of more hidden layers or neurons in each hidden layer to the chosen configuration was found to yield similar MSE errors or accuracies and hence are not considered in this study.

## 5. Results

Followed by the determination of FFNN configuration, the rank of the features in each dataset is evaluated using the proposed method. Furthermore, other feature ranking methods are also considered in this study for the sake of comparison. An open-source software WEKA is employed for this purpose. While feature ranking methods such as Pearson correlation coefficient, ReliefF and, mutual information are used for regression task, symmetric uncertainty, information gain, gain ratio, reliefF and, chi-square is employed for the classification task. The efficacy of all feature ranking methods is then assessed by evaluating the performance of FFNN, wherein the size of the input layer is increased by one feature in each succession. In other words, the performance of FFNN for the only top-most feature is first assessed, and then the process is repeated by including the second most important feature and so on.

### 5.1 Regression

From Table 3, it can be inferred that all four feature ranking methods yielded feature 6 (Abdomen) as the most important feature and feature 10 (Ankle) as the least relevant feature for determining the percentage of body fat. While the top six features determined using Pearson correlation coefficient, ReliefF and, mutual information method are noticed to be similar; the proposed method yielded different feature ranks. Furthermore, the MSE for body fat dataset with each feature's inclusion is evaluated for all four feature ranking methods and is shown in Figure 3(a). From Figure 3(a), it is evident that the overall trend of MSE for FFNN decreases with the inclusion of each feature. While the proposed method was found to yield lower MSE with only seven top-most features, the mutual information method yielded lower MSE for eleven features for the bodyfat dataset. In other words, the filter based approach was found to be ineffective at determining a subset of important features that could reduce the MSE. According to the proposed method, following features are found to be least important as they do not contribute further for reduction of MSE: (5) Chest (cm), (7) Hip (cm), (9) Knee (cm), (10) Ankle (cm), (11) Biceps (cm), (12) Forearm (cm).

In the case of the abalone dataset, the least relevant features are determined to be the same by all four feature ranking methods, i.e., feature 1 (female), feature 2 (infant), and feature 3 (male) are identified to be the least relevant (see Table 3). While the remaining seven features' rank was found to vary, feature 10 (shell weight) and feature 7 (whole weight) were common in the top four features for all feature ranking methods, including the proposed method. Similar to the body fat dataset, the MSE of FFNN with the inclusion of each feature is determined for all feature ranking methods and is shown in Figure 3(b). From Figure 3 (b), it can be inferred that the trend of ReliefF and the proposed method are similar. Both ReliefF and the proposed method identified feature 5



(diameter), feature 6 (height), feature 7 (whole weight), and feature 10 (shell weight) as the top 4 features that yield the lowest MSE. In other words, ReliefF was found to be effective among all the filter-based methods. According to the proposed method, following features are found to be least important as they do not contribute further for reduction of MSE: (1) Female, (2) Infant, (3) Male, (4) Length (gms.), (8) Shucked weight (gms.), (9) Viscera weight (gms.).

Interestingly, in the wine quality dataset, all four feature ranking methods yielded different ranks for the features (see Table 3). However, feature 11 (alcohol) is determined to be one of the top two features by all four feature ranking methods. Furthermore, feature 6 (free sulfurdioxide) is determined to be common among first four features determined by all feature ranking methods except mutual information. The MSE of FFNN with each feature's inclusion is determined for all feature ranking methods and is shown in Figure 3(c). The trend obtained in Figure 3 (c), reveals that all feature ranking methods performed more or less similar.

## 5.2 Classification

From Table 4, it can be inferred that all feature ranking methods employed for the classification task identified similar least relevant features for the vehicle dataset (feature 15 (skewness minor), feature 16 (kurtosis major)). However, the rank of the remaining features was found to vary. While feature 12 (scaled variance minor), feature 7 (scatter ratio) and feature 8 (elongatedness) was found to be the top three features for symmetric uncertainty, information gain, gain ratio, reliefF and, chi-square, feature 10 (maximum length rectangularity), feature 8 (elongatedness) and feature 5 (axis aspect ratio) was found to be the top 3 features for the proposed method, i.e., feature 8 (elongatedness) was found to be common among top 3 features predicted by all feature ranking methods. Furthermore, the trend of the accuracy is determined for vehicle dataset for all feature ranking methods with the inclusion of each feature in succession and is shown in Figure 4(a). From Figure 4(a), it is evident that the accuracy of the FFNN increases with the addition of each feature for the vehicle dataset. The proposed method yielded an accuracy of 75% by selecting only the top 6 features and was found to outperform the other feature ranking methods. The top 6 features are identified as follows: (5) axis aspect ratio, (8) elongatedness, (10) maximum length rectangularity, (14) skewness major, (17) kurtosis minor and (18) hollow ratio.

Similar to the vehicle dataset, all feature ranking methods employed in the case of the segmentation dataset obtained the same least relevant features (feature 1 (region-centroid-col), feature 3 (short-line density), feature 4 (lines of length), feature 6 (vedge-sd), and feature 8 (hedge-sd)). While the rank of the top features was found to vary for all feature ranking methods, feature 10 (rawred-mean), feature 16 (value-mean), and feature 18 (hue-mean) were found to be common among the top-most 6 features. The trend of the accuracy for the segmentation dataset is determined for all feature ranking methods with the inclusion of each feature in succession and is shown in Figure 4(b). From Figure 4(b), it is evident that the accuracy of the FFNN increases with the addition of each feature for the segmentation dataset. Among all the feature ranking methods, the proposed method was found to outperform yielding the highest accuracy of 90% with only the top 6 features. In other words, the filter based methods suggested top 10 features are important for achieving an accuracy of 85%.

Interestingly, in the breast cancer dataset, all feature ranking methods resulted in similar top-most features, i.e., feature 21 (radius3) and feature 23 (perimeter3). While symmetric uncertainty, information gain, gain ratio, reliefF and, chi-square identified feature 10 (fractal dimension1), feature 12 (texture2), and feature 15 (smoothness2) as least relevant, the proposed method identified the feature 3 (perimeter1), feature 5 (smoothness1) and feature 27 (concavity3) are least relevant. Similar to the vehicle and segmentation dataset, the trend of accuracy is obtained for the breast cancer dataset with the inclusion of each feature in each succession and is shown in Figure 4(c). In the case of the breast cancer dataset, the trend of all feature ranking methods was found to be more or less similar. An accuracy of 93% is achieved by the inclusion of the top two features, i.e., feature 21 (radius3) and feature 23 (perimeter3).

## 6. Summary and Future Work

A novel complex-step sensitivity analysis-based feature selection method is proposed in this study for regression and classification tasks. A step-by-step process involved in implementing the proposed method in the framework of FFNN is described, and its efficacy on real-world datasets is demonstrated. Three real-world datasets, namely, body fat percentage dataset, abalone dataset, and wine quality dataset, are chosen for the regression task and, three datasets, namely vehicle dataset, segmentation dataset, and breast cancer dataset, are chosen for the classification task. While the proposed method was found to outperform other popular feature ranking methods for classification datasets (vehicle, segmentation, and breast cancer), it was found to perform more or less similar with other methods in the case of regression datasets (body fat, abalone, and wine quality). An average MSE of 20.41, 4.6, and 0.53 were observed for body fat, abalone, and wine quality datasets, respectively, and an average accuracy of 75%, 80%, and 90% was observed for the vehicle segmentation and breast cancer datasets, respectively. Furthermore, the top-most relevant features and irrelevant features are identified for all the employed datasets. At this juncture, it is also important to note that the proposed method possesses the advantage of performing sensitivity analysis through the forward propagation of FFNN, i.e., no backpropagation is required for evaluating the derivatives.

In future work, the authors intend to extend the proposed method to the multiple output regression problems. In addition to this, the authors would also like to investigate the influence of different activation functions (e.g., Sigmoid, tanh, Softplus, Leaky ReLU, etc.). Other supervised ML classification algorithms will be employed, and the efficacy of the proposed method will be examined. Note that often complete dataset may not be required for training the FFNN when the size of the dataset is large. Hence the influence of a number of instances on the determination of the important features would also be studied. Furthermore, the proposed method would also be extended to the datasets that consists of discrete and continuous features and also include redundant features.

## Data Availability Statement

All the datasets employed in this study are obtained from UCI open-source data repositories [49–53].

## Acknowledgment

337 Research presented in this paper was supported by the National Science Foundation under NSF  
338 EPSCoR Track-1 Cooperative Agreement OIA #1946202. Any opinions, findings, and  
339 conclusions, or recommendations expressed in this material are those of the author(s) and do not  
340 necessarily reflect the views of the National Science Foundation.

341

## References

- [1] Dash M, Liu H. Feature Selection for Classification. vol. 1. 1997.
- [2] Blum AL, Langley P. Selection of Relevant Features and Examples in Machine Learning. n.d.
- [3] Sánchez-Marño N, Alonso-Betanzos A, Tombilla-Sanromán M. Filter methods for feature selection - A comparative study. Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 4881 LNCS, Springer Verlag; 2007, p. 178–87. [https://doi.org/10.1007/978-3-540-77226-2\\_19](https://doi.org/10.1007/978-3-540-77226-2_19).
- [4] Chandrashekar G, Sahin F. A survey on feature selection methods q. Comput Electr Eng 2014;40:16–28. <https://doi.org/10.1016/j.compeleceng.2013.11.024>.
- [5] Bolón-Canedo V, Sánchez-Marño N, Alonso-Betanzos A. A review of feature selection methods on synthetic data. Knowl Inf Syst 2013;34:483–519. <https://doi.org/10.1007/s10115-012-0487-8>.
- [6] Tadist K, Nikolov NS, Mrabti F, Zahi A. Feature selection methods and genomic big data: a systematic review n.d. <https://doi.org/10.1186/s40537-019-0241-0>.
- [7] Perscheid C, Grasnick B, Uflacker M. Integrative Gene Selection on Gene Expression Data: Providing Biological Context to Traditional Approaches. J Integr Bioinform 2018;16. <https://doi.org/10.1515/jib-2018-0064>.
- [8] Jović A, Brkić K, Bogunović N. A review of feature selection methods with applications. 2015 38th Int. Conv. Inf. Commun. Technol. Electron. Microelectron. MIPRO 2015 - Proc., Institute of Electrical and Electronics Engineers Inc.; 2015, p. 1200–5. <https://doi.org/10.1109/MIPRO.2015.7160458>.
- [9] Asir Antony Gnana Singh D, Appavu alias Balamurugan SK, Jebamalar Leavline E. Literature Review on Feature Selection Methods for High-Dimensional Data. vol. 136. 2016.
- [10] Naik DL, Sajid HU, Kiran R. Texture-Based Metallurgical Phase Identification in Structural Steels: A Supervised Machine Learning Approach. Metals (Basel) 2019;9:546. <https://doi.org/10.3390/met9050546>.
- [11] Naik DL, Kiran R. Naïve Bayes classifier, multivariate linear regression and experimental testing for classification and characterization of wheat straw based on mechanical properties. Ind Crops Prod 2018;112. <https://doi.org/10.1016/j.indcrop.2017.12.034>.
- [12] Naik DL, Kiran R. Identification and characterization of fracture in metals using machine learning based texture recognition algorithms. Eng Fract Mech 2019;219. <https://doi.org/10.1016/j.engfracmech.2019.106618>.
- [13] Dong NT, Winkler L, Khosla M. Revisiting Feature Selection with Data Complexity for Biomedicine n.d. <https://doi.org/10.1101/754630>.
- [14] Hua J, Tembe WD, Dougherty ER. Performance of feature-selection methods in the classification of high-dimension data. Pattern Recognit 2009;42:409–24.

380 <https://doi.org/10.1016/j.patcog.2008.08.001>.

381 [15] Cilia N, De Stefano C, Fontanella F, Raimondo S, Scotto di Freca A. An Experimental  
382 Comparison of Feature-Selection and Classification Methods for Microarray Datasets.  
383 Information 2019;10:109. <https://doi.org/10.3390/info10030109>.

384 [16] Hira ZM, Gillies DF. A Review of Feature Selection and Feature Extraction Methods  
385 Applied on Microarray Data 2015. <https://doi.org/10.1155/2015/198363>.

386 [17] Remeseiro B, Bolon-Canedo V. A review of feature selection methods in medical  
387 applications. Comput Biol Med 2019;112:103375.  
388 <https://doi.org/10.1016/j.combiomed.2019.103375>.

389 [18] Garrett D, Peterson DA, Anderson CW, Thaut MH. Comparison of linear, nonlinear, and  
390 feature selection methods for EEG signal classification. IEEE Trans Neural Syst Rehabil  
391 Eng 2003;11:141–4. <https://doi.org/10.1109/TNSRE.2003.814441>.

392 [19] Refaeilzadeh P, Tang L, Liu H. On Comparison of Feature Selection Algorithms. n.d.

393 [20] Liu J, Wang G. A hybrid feature selection method for data sets of thousands of variables.  
394 Proc. - 2nd IEEE Int. Conf. Adv. Comput. Control. ICACC 2010, vol. 2, 2010, p. 288–91.  
395 <https://doi.org/10.1109/ICACC.2010.5486671>.

396 [21] Hoque N, Singh M, Bhattacharyya DK. EFS-MI: an ensemble feature selection method  
397 for classification. Complex Intell Syst 2018;4:105–18. [https://doi.org/10.1007/s40747-](https://doi.org/10.1007/s40747-017-0060-x)  
398 017-0060-x.

399 [22] Setiono R, Liu H. Neural-network feature selector. IEEE Trans Neural Networks  
400 1997;8:654–62. <https://doi.org/10.1109/72.572104>.

401 [23] Sindhwani V, Rakshit S, Deodhare D, Erdogmus D, Principe JC, Niyogi P. Feature  
402 selection in MLPs and SVMs based on maximum output information. IEEE Trans Neural  
403 Networks 2004;15:937–48. <https://doi.org/10.1109/TNN.2004.828772>.

404 [24] Bo L, Wang L, Jiao L. Multi-layer perceptrons with embedded feature selection with  
405 application in cancer classification. Chinese J Electron 2006;15:832–5.

406 [25] Ruck DW, Rogers SK, Kabrisky M. Feature Selection Using a Multilayer Perceptron. vol.  
407 2. 1990.

408 [26] Gasca E, Sánchez JS, Alonso R. Eliminating redundancy and irrelevance using a new  
409 MLP-based feature selection method. Pattern Recognit 2006;39:313–5.  
410 <https://doi.org/10.1016/j.patcog.2005.09.002>.

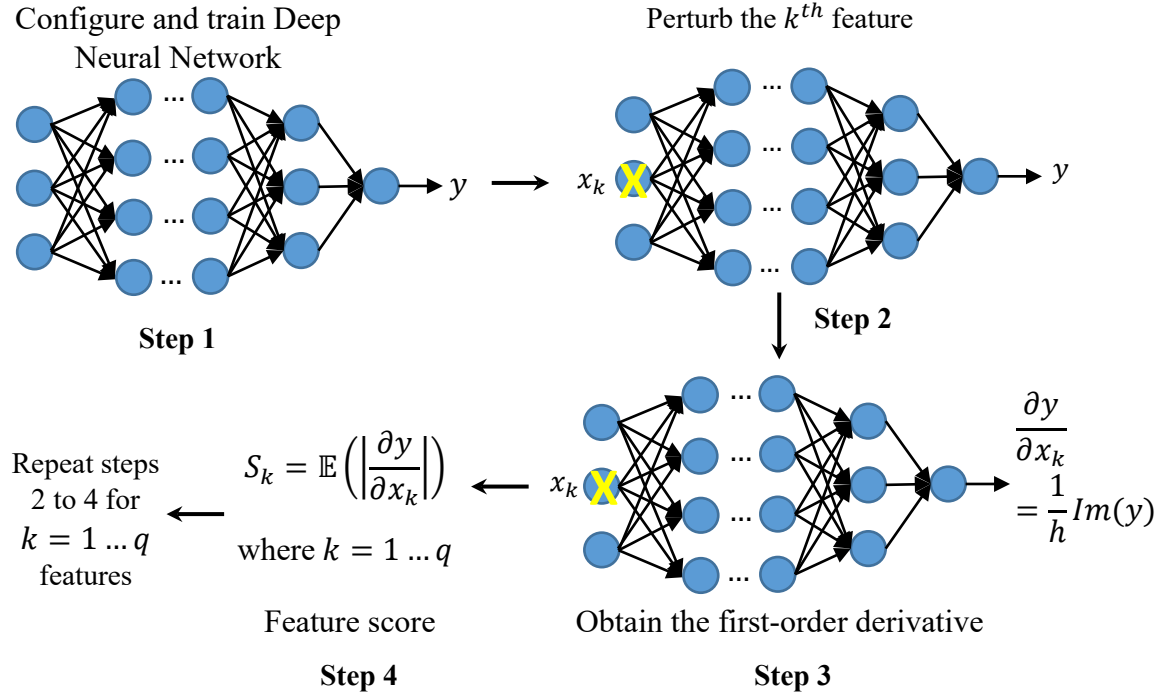
411 [27] Utans J, Moody J, Rehfuess S, Siegelmann H. Input variable selection for neural networks:  
412 application to predicting the U.S. business cycle. IEEE/IAFE Conf. Comput. Intell.  
413 Financ. Eng. Proc., IEEE; 1995, p. 118–22. <https://doi.org/10.1109/cifer.1995.495263>.

414 [28] Hajnayeb A, Ghasemloonia A, Khadem SE, Moradi MH. Application and comparison of  
415 an ANN-based feature selection method and the genetic algorithm in gearbox fault  
416 diagnosis. Expert Syst Appl 2011;38:10205–9.  
417 <https://doi.org/10.1016/j.eswa.2011.02.065>.

- 418 [29] Acir N. A support vector machine classifier algorithm based on a perturbation method and  
 419 its application to ECG beat recognition systems. *Expert Syst Appl* 2006;31:150–8.  
 420 <https://doi.org/10.1016/J.ESWA.2005.09.013>.
- 421 [30] Montaña JJ, Palmer A. Numeric sensitivity analysis applied to feedforward neural  
 422 networks. *Neural Comput Appl* 2003;12:119–25. [https://doi.org/10.1007/s00521-003-](https://doi.org/10.1007/s00521-003-0377-9)  
 423 [0377-9](https://doi.org/10.1007/s00521-003-0377-9).
- 424 [31] Güne, A, Baydin G, Pearlmutter BA, Siskind JM. *Automatic Differentiation in Machine*  
 425 *Learning: a Survey*. 2018.
- 426 [32] Jerrell ME. Automatic Differentiation and Interval Arithmetic for Estimation of  
 427 Disequilibrium Models. *Comput Econ* 1997;10:295–316.  
 428 <https://doi.org/10.1023/A:1008633613243>.
- 429 [33] Hashem S. Sensitivity analysis for feedforward artificial neural networks with  
 430 differentiable activation functions, *Institute of Electrical and Electronics Engineers*  
 431 *(IEEE)*; 2003, p. 419–24. <https://doi.org/10.1109/ijcnn.1992.287175>.
- 432 [34] Driscoll TA, Braun RJ. *Fundamentals of Numerical Computation*. 2017.
- 433 [35] Boudjemaa R, Cox MG, Forbes AB, Harris PM. Report to the National Measurement  
 434 Directorate, Department of Trade and Industry From the Software Support for Metrology  
 435 Programme Automatic Differentiation Techniques and their Application in Metrology.  
 436 2003.
- 437 [36] CANUL-REICH J, HALL LO, GOLDFOF DB, KORECKI JN, ESCHRICH S.  
 438 ITERATIVE FEATURE PERTURBATION AS A GENE SELECTOR FOR  
 439 MICROARRAY DATA. [Http://DxDoiOrg/101142/S0218001412600038](http://DxDoiOrg/101142/S0218001412600038) 2012;26.  
 440 <https://doi.org/10.1142/S0218001412600038>.
- 441 [37] Lyness JN, Moler CB. Numerical Differentiation of Analytic Functions. *SIAM J Numer*  
 442 *Anal* 1967;4:202–10. <https://doi.org/10.1137/0704019>.
- 443 [38] Squire W, Trapp G. Using complex variables to estimate derivatives of real functions.  
 444 *SIAM Rev* 1998;40:110–2. <https://doi.org/10.1137/S003614459631241X>.
- 445 [39] Kiran R, Khandelwal K. Complex step derivative approximation for numerical evaluation  
 446 of tangent moduli. *Comput Struct* 2014;140:1–13.  
 447 <https://doi.org/10.1016/j.compstruc.2014.04.009>.
- 448 [40] Kiran R, Khandelwal K. Automatic implementation of finite strain anisotropic  
 449 hyperelastic models using hyper-dual numbers. *Comput Mech* 2015;55:229–48.  
 450 <https://doi.org/10.1007/s00466-014-1094-1>.
- 451 [41] Martins J, Sturdza P, Alonso J, R A Martins JR, Alonso JJ. The complex-step derivative  
 452 approximation. *ACM Trans Math Software, Assoc Comput Mach* 2003;29:245–62.  
 453 <https://doi.org/10.1145/838250.838251i>.
- 454 [42] Conolly J, Lake M. *Geographical information systems in archaeology* 2006:338.
- 455 [43] Zhu J-J, Li J-J. *Quantum Dots*, 2013. [https://doi.org/10.1007/978-3-642-44910-9\\_2](https://doi.org/10.1007/978-3-642-44910-9_2).

- [44] Campbell AR. Numerical Analysis of Complex-Step Differentiation in Spacecraft Trajectory Optimization Problems. 2011.
- [45] Kiran R, Li L, Khandelwal K. Complex Perturbation Method for Sensitivity Analysis of Nonlinear Trusses. *J Struct Eng* 2017;143:04016154. [https://doi.org/10.1061/\(asce\)st.1943-541x.0001619](https://doi.org/10.1061/(asce)st.1943-541x.0001619).
- [46] Lai KL, Crassidis JL, Cheng Y, Kim J. New complex-step derivative approximations with application to second-order Kalman filtering. *Collect. Tech. Pap. - AIAA Guid. Navig. Control Conf.*, vol. 2, 2005, p. 982–98. <https://doi.org/10.2514/6.2005-5944>.
- [47] Christopher MB. *Neural Networks for Pattern Recognition*. Oxford University Press; 1995.
- [48] UCI Machine Learning Repository n.d. <https://archive.ics.uci.edu/ml/index.php> (accessed April 7, 2021).
- [49] Johnson RW. Fitting Percentage of Body Fat to Simple Body Measurements. *J Stat Educ* 1996;4. <https://doi.org/10.1080/10691898.1996.11910505>.
- [50] Nash WJ. The Population Biology of Abalone (*Haliotis* species) in Tasmania. I. Blacklip Abalone (*H. rubra*) from the North Coast and Islands of Bass Strait. 1994.
- [51] Cortez P, Cerdeira A, Almeida F, Matos T, Reis J. Modeling wine preferences by data mining from physicochemical properties. *Decis Support Syst* 2009;47:547–53. <https://doi.org/10.1016/j.dss.2009.05.016>.
- [52] Siebert JP. *Vehicle Recognition Using Rule Based Methods* 1987.
- [53] Mangasarian OL, Street WN, Wolberg WH. Breast Cancer Diagnosis and Prognosis Via Linear Programming. *Oper Res* 1995;43:570–7. <https://doi.org/10.1287/opre.43.4.570>.
- [54] Ding B, Qian H, Zhou J. Activation functions and their characteristics in deep neural networks. *Proc. 30th Chinese Control Decis. Conf. CCDC 2018*, Institute of Electrical and Electronics Engineers Inc.; 2018, p. 1836–41. <https://doi.org/10.1109/CCDC.2018.8407425>.

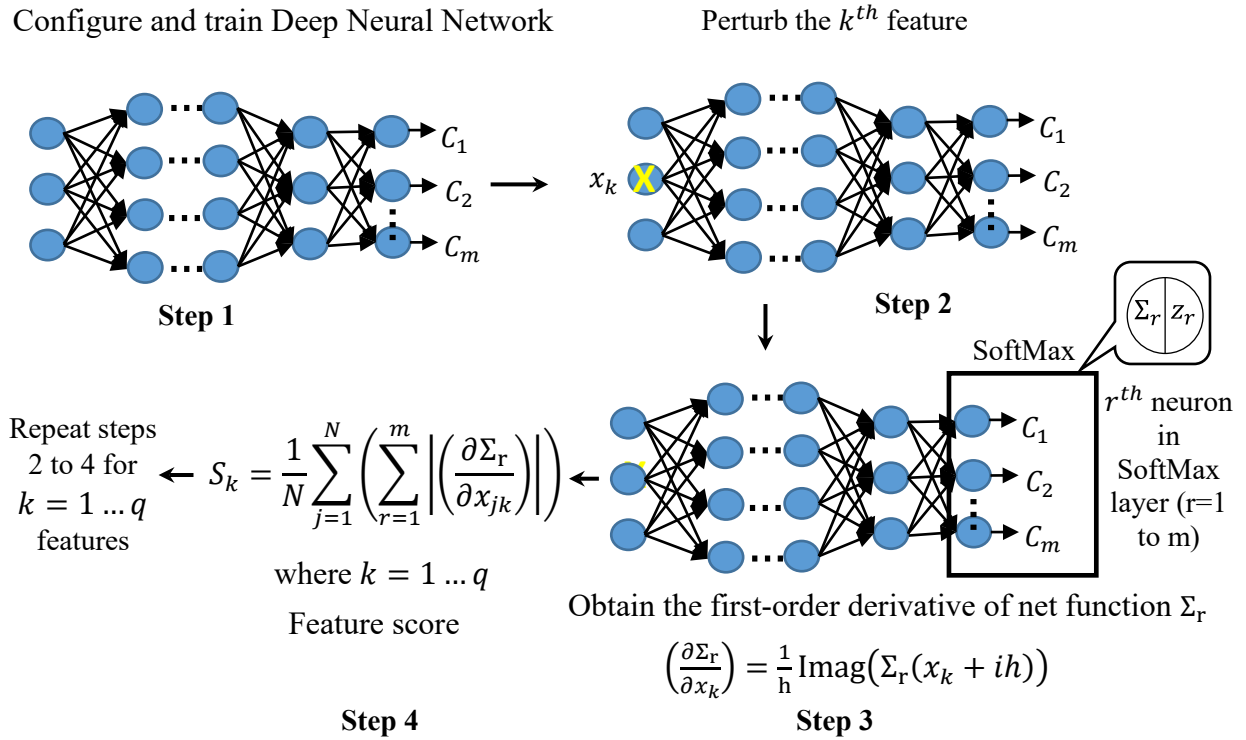
483 **Figures**



484

485

Figure 1. Steps involved in the complex-step sensitivity for regression task.



486

487

Figure 2. Steps involved in the complex-step sensitivity for the classification task.



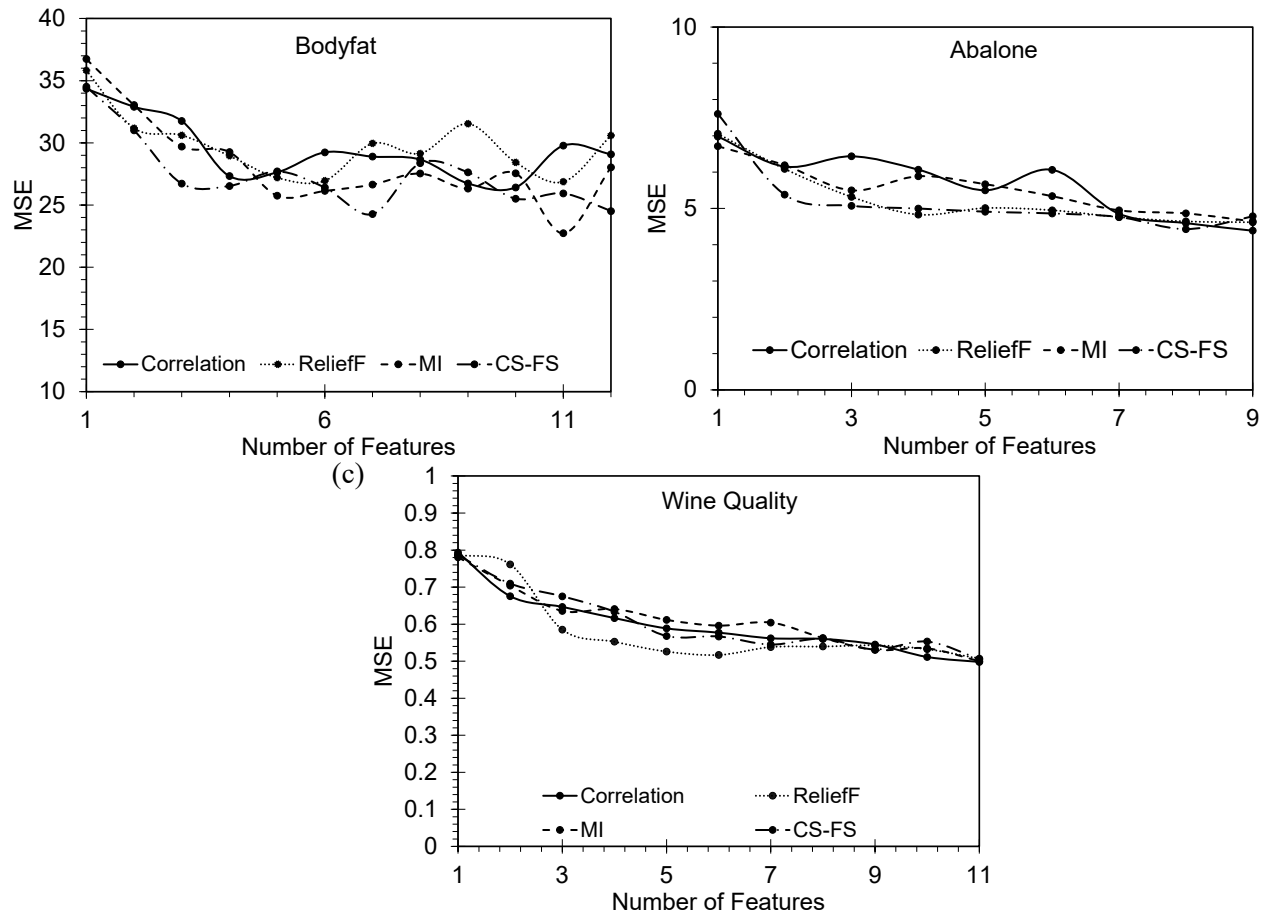


Figure 3. Comparison of the complex-step sensitivity method with other feature selection methods for regression task.

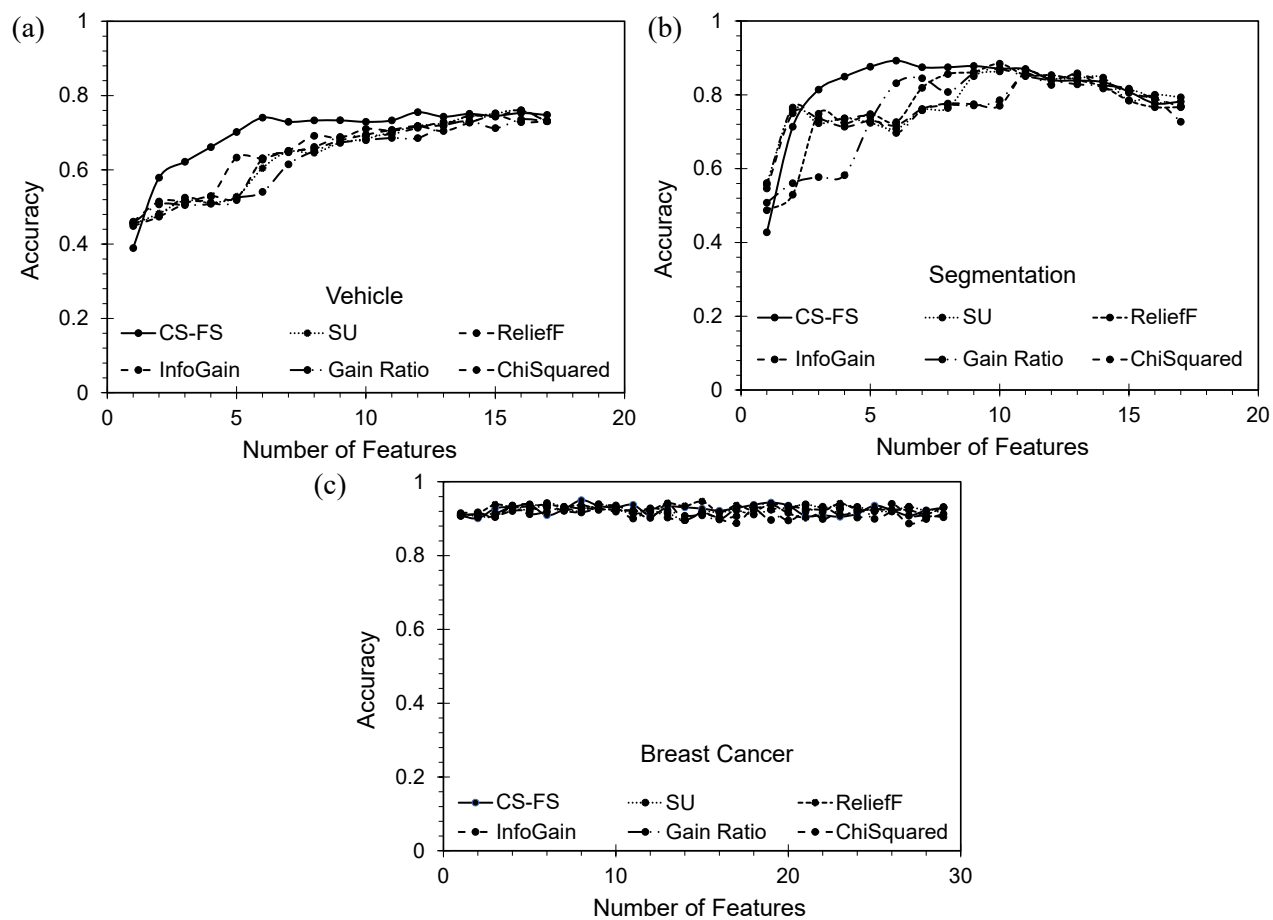


Figure 4. Comparison of the complex-step sensitivity method with other feature selection methods for the classification task.

## Tables

Table 1. Description of the datasets used for regression task.

Dataset name	Instances	No. of features	No. of target variables
Bodyfat	252	13	1
Abalone	4177	10	1
Wine quality	4898	11	1

Table 2. Description of the datasets used for the classification task.

Dataset name	Instances	No. of features	No. of class labels
Vehicle	846	30	4
Segmentation	210	18	7
Breast cancer	569	18	2

Table 3. Important features identified by various feature selection methods for regression task (ranked in the descending order of their importance).

Bodyfat dataset				Abalone dataset				Wine quality dataset			
Corr.	ReliefF	MI	CS-FS	Corr.	ReliefF	MI	CSDA	Corr.	ReliefF	MI	CS-FS
6	6	6	6	10	10	10	7	11	2	8	11
5	5	5	3	5	7	5	8	9	11	11	4
7	7	7	13	6	8	7	6	10	6	4	6
2	2	2	4	4	9	6	10	6	9	7	2
8	8	8	8	7	5	9	9	3	7	5	7
9	9	9	2	9	6	4	4	4	1	6	5
1	1	11	1	8	4	8	5	1	10	3	1
11	11	4	7	1	2	2	2	7	8	2	9
3	3	1	5	3	3	1	3	2	3	9	3
4	4	13	12	2	1	3	1	5	4	1	8
10	10	12	11					8	5	10	10
13	13	10	10								
12	12	3	9								

Table 4. Important features identified by various feature selection methods for classification task (ranked in the descending order of their importance).

	Method	Feature Ranking
Vehicle dataset	ReliefF	8, 7, 12, 9, 3, 11, 18, 4, 2, 1, 13, 10, 16, 14, 17, 6, 15, 5
	Symmetric Uncertainty	12, 7, 8, 11, 9, 6, 3, 4, 1, 13, 2, 14, 10, 17, 18, 5, 16, 15
	Info Gain	12, 7, 8, 11, 9, 3, 6, 2, 1, 4, 13, 10, 14, 17, 18, 5, 16, 15
	Gain Ratio	11, 9, 12, 7, 4, 8, 6, 3, 5, 18, 13, 14, 1, 2, 16, 10, 15, 17
	Chi-Squared	12, 7, 8, 9, 11, 3, 6, 1, 2, 10, 14, 13, 4, 17, 18, 5, 16, 15
	CSDA	10, 8, 5, 17, 14, 18, 11, 3, 6, 12, 7, 1, 9, 4, 13, 2, 15, 16
Segmentation dataset	ReliefF	11, 16, 18, 9, 12, 10, 2, 15, 14, 13, 17, 1, 5, 7, 3, 4, 6, 8
	Symmetric Uncertainty	18, 10, 9, 16, 12, 11, 15, 17, 2, 14, 13, 7, 8, 5, 6, 3, 4, 1
	Info Gain	18, 9, 12, 16, 10, 11, 15, 17, 13, 14, 2, 7, 8, 5, 6, 3, 4, 1
	Gain Ratio	10, 11, 9, 16, 18, 2, 12, 14, 15, 17, 13, 8, 7, 5, 6, 3, 4, 1
	Chi-Squared	18, 12, 9, 16, 10, 11, 13, 15, 17, 14, 2, 7, 8, 5, 6, 3, 4, 1
	CSDA	2, 18, 15, 13, 10, 16, 11, 12, 17, 9, 14, 6, 8, 7, 5, 4, 3, 1
Breast cancer dataset	ReliefF	28, 8, 21, 23, 3, 1, 7, 24, 4, 27, 26, 6, 22, 25, 11, 2, 14, 13, 29, 30, 10, 18, 5, 16, 9, 17, 19, 15, 12, 20.
	Symmetric Uncertainty	23, 21, 24, 28, 8, 3, 7, 4, 1, 27, 14, 11, 13, 6, 26, 17, 2, 18, 22, 25, 29, 16, 5, 30, 9, 19, 20, 10, 12, 15.
	Info Gain	23, 24, 21, 28, 8, 3, 4, 1, 7, 14, 27, 11, 13, 26, 6, 17, 18, 22, 2, 29, 16, 25, 9, 5, 30, 20, 19, 10, 12, 15.
	Gain Ratio	23, 21, 24, 28, 8, 7, 27, 3, 4, 1, 14, 6, 11, 13, 26, 17, 2, 19, 18, 25, 22, 29, 5, 16, 30, 9, 20, 12, 10, 15.
	Chi-Squared	23, 21, 24, 28, 8, 3, 4, 1, 7, 14, 27, 11, 13, 26, 6, 17, 18, 22, 2, 29, 25, 16, 9, 5, 30, 20, 19, 10, 12, 15.
	CSDA	21, 23, 28, 20, 8, 4, 7, 11, 24, 17, 15, 2, 22, 30, 12, 26, 13, 16, 1, 14, 10, 9, 29, 25, 18, 19, 6, 3, 27, 5.