

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

# SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification

M.Z. Naser<sup>1,2</sup>, Mohammad Khaled al-Bashiti<sup>1</sup>, A.Z. Naser<sup>3,4</sup>

<sup>1</sup>School of Civil & Environmental Engineering and Earth Sciences (SCEEES), Clemson University, USA

<sup>2</sup>Artificial Intelligence Research Institute for Science and Engineering (AIRISE), Clemson University, USA

E-mail: [mznaser@clemson.edu](mailto:mznaser@clemson.edu), [malbash@g.clemson.edu](mailto:malbash@g.clemson.edu), Website: [www.mznaser.com](http://www.mznaser.com)

<sup>3</sup>School of Engineering, University of Guelph, Guelph, ON N1G 2W1, Canada, E-mail: [anaser@uoguelph.ca](mailto:anaser@uoguelph.ca)

<sup>4</sup>Department of Mechanical Engineering, University of Manitoba, Winnipeg R3T 5V6, Canada

## Abstract

The field of machine learning (ML) has witnessed significant advancements in recent years. However, many existing algorithms lack interpretability and struggle with high-dimensional and imbalanced data. This paper proposes SPINEX, a novel similarity-based interpretable neighbor exploration algorithm designed to address these limitations. This algorithm combines ensemble learning and feature interaction analysis to achieve accurate predictions and meaningful insights by quantifying each feature's contribution to predictions and identifying interactions between features, thereby enhancing the interpretability of the algorithm. To evaluate the performance of SPINEX, extensive experiments on 59 synthetic and real datasets were conducted for both regression and classification tasks. The results demonstrate that SPINEX achieves comparative performance and, in some scenarios, may outperform commonly adopted ML algorithms. The same findings demonstrate the effectiveness and competitiveness of SPINEX, making it a promising approach for various real-world applications.

**Keywords:** Algorithm; Machine learning; Interpretability; Supervised learning.

## 1.0 Introduction

The rapid growth of machine learning (ML) techniques has revolutionized various domains, enabling accurate predictions and decision-making [1,2]. However, challenges persist, such as the lack of interpretability in complex models [3]. This has led to a growing interest in interpretable ML algorithms [4]. While existing approaches, such as decision trees and linear models, offer interpretability, they often sacrifice predictive performance. Conversely, complex models like neural networks and ensemble methods achieve high accuracy but lack interpretability [5].

In the vast landscape of algorithmic design, similarity-based algorithms are potent methods for tackling various prediction problems [6,7]. These algorithms rely on the premise that similar objects share similar properties and hence draw their strength from leveraging instance similarities and proximity to make predictions or categorizations/clustering. This approach is particularly useful in recommendation systems, classification tasks, and regression problems, where the goal is to predict an outcome based on the similarity of the input data to previously seen examples.

The core idea behind such models is to find the 'neighbors' of a given data point in the feature space. Once these neighbors are identified, they are used to predict the given data point. This can be realized by taking a weighted average of the neighbors' outcomes, where the weights are

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

39 determined by the similarity of each neighbor to the given data point. This underlying principle is  
40 harnessed in various contexts, such as document retrieval [8], image recognition [9], etc.

41 One of the key advantages of these models is their explainability. Unlike many other ML models,  
42 similarity-based models can provide clear explanations for their predictions based on specific,  
43 identifiable neighbors rather than abstract mathematical functions. However, the concept of  
44 'similarity' is not always straightforward. Different similarity measures may be appropriate for  
45 different data types, and choosing the right measure (i.e., including the number of neighbors) can  
46 significantly impact the model's performance [10].

47 In addition to these challenges, these models also face the 'curse of dimensionality'. This refers to  
48 the fact that as the number of features (dimensions) in the data increases, similarity can become  
49 less meaningful. This is because, in high-dimensional spaces, all data points tend to be 'far away'  
50 from each other, making it difficult to find meaningful neighbors. Despite these challenges, recent  
51 research has shown that these models can achieve superior performance on a variety of tasks. For  
52 example, these models have been successfully applied to recommendation systems and  
53 classification tasks, where they have been shown to outperform traditional collaborative filtering  
54 approaches [11,12]. Notable reviews on this front can be found elsewhere [13–15].

55 As mentioned above, one way to implement the similarity-based approach is to use a nearest-  
56 neighbor algorithm. The nearest neighbor algorithm finds the  $k$  most similar objects to a given  
57 object and then predicts the label of the given object based on the labels of the  $k$  nearest neighbors.  
58 Such algorithms rely on the idea that the characteristics of an object can be inferred from those of  
59 its close neighbors [16]. In a  $kNN$  algorithm, the class of a new object is determined by the classes  
60 of its  $k$  nearest neighbors in the training set. The "closeness" of objects is typically defined in terms  
61 of some distance measure, such as Euclidean distance for numeric data or Hamming distance for  
62 categorical data [17]. The strength of such methods is their simplicity and intuitiveness since their  
63 predictions are directly tied to the observed data [18,19].

64 On a different front, ML based on similarity computation approaches was reported to achieve  
65 promising experimental results in Natural Language Processing (NLP). Mikolov et al. [32]  
66 introduced two models that were designed to calculate continuous vector representations of words  
67 from extensive datasets via similarity tasks and showed superior accuracy to various neural  
68 networks. Also, the new similarity-based models were reported to outperform former models in  
69 terms of computational cost as they require less than a day to learn high quality word vectors from  
70 a more than 1.5 billion words dataset, demonstrating high merit in utilizing semantic word  
71 similarities.

## 72 *1.1 Related works*

73 Central to similarity-based algorithms' functionality is the concept of identifying the nearest  
74 neighbors to a given data point of interest; these neighbors collectively contribute to the prediction  
75 [20]. This central characteristic is particularly advantageous in scenarios where the underlying data

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

76 distribution is unknown or complex and allows direct application in pattern recognition [21],  
77 recommendation systems [22], and medical diagnosis [23,24].

78 Over the years, various studies have explored the effectiveness of similarity-based algorithms. One  
79 significant area of research has been related to weighted nearest neighbors, where different  
80 neighbors are assigned different weights based on their distance from the target point [25]. This  
81 approach mitigates some of the limitations of the basic-like algorithms, which treat all neighbors  
82 equally (regardless of their proximity to the target point). Another aspect has been the exploration  
83 of efficient algorithms for large datasets. This stems from the fact that the brute-force approach of  
84 comparing a target point with every other point in a large dataset can be computationally  
85 impractical [26]. This challenge has led to the development of various indexing structures and  
86 algorithms [27].

87 The integration of similarity-based algorithms with other ML techniques, such as deep learning,  
88 presents an exciting theme for exploration [28,29]. For example, recent works have adopted  
89 similarity-based algorithms with clustering [30] and anomaly detection [31]. Dudek and Pelka [7]  
90 explored the application of pattern similarity-based models. They noted that experimental results  
91 collected across 35 European countries showed that such models outperformed the classical  
92 statistical and ML models in terms of accuracy, simplicity, and ease of optimization.

### 93 *1.2 Key contributions of this work*

94 From the lens of this paper, we propose a novel ML algorithm, SPINEX (Similarity-based  
95 Predictions with Explainable Neighbors Exploration), which combines similarity-based predictions  
96 and neighbors' exploration. SPINEX offers interpretability through feature contribution analysis and  
97 interaction effects. A number of extensive experiments were carried out to validate the  
98 effectiveness and competitiveness of SPINEX in both regression and classification tasks. More  
99 specifically, to evaluate the performance of SPINEX, 59 experiments were conducted on diverse  
100 synthetic and real datasets covering a wide range of domains. The experimental results  
101 demonstrate the effectiveness and competitiveness of SPINEX compared to state-of-the-art  
102 algorithms. Therefore, this paper contributes to the field by proposing a novel ML algorithm that  
103 combines the strengths of similarity-based predictions and neighbors' exploration, offers  
104 interpretability, and demonstrates its effectiveness and competitiveness in regression and  
105 classification tasks.

106 The rest of the paper is organized as follows: Section 2 presents the methodology of SPINEX,  
107 explaining each component in detail. Section 3 discusses the experimental setup and presents the  
108 results of the experiments conducted. Finally, Section 4 concludes the paper, highlighting the  
109 contributions of SPINEX and suggesting potential future directions for research.

## 110 **2.0 Description of SPINEX**

111 The SPINEX algorithm comprises several components that provide interpretable regression and  
112 classification analysis. First, it begins with a data preprocessing step, then calculates pairwise

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

113 distances between instances using a user-defined metric. Based on these distances, weights are  
114 assigned to the instances using a Gaussian kernel function. This step allows SPINEX to emphasize  
115 the influence of the most relevant neighbors during prediction. SPINEX also accommodates single  
116 and ensemble models to capture the inherent complexity of the data. This algorithm builds on the  
117 concept of neighbor-based feature importance, which measures the contribution of each feature to  
118 the prediction by considering the influence of its neighboring instances [33]. This approach  
119 provides a more nuanced understanding of feature importance, accounting for individual feature  
120 effects and their dependencies on nearby instances.

121 SPINEX also incorporates feature interaction analysis, which explores the interactions between  
122 different feature combinations to identify synergistic or antagonistic effects on the target variable.  
123 Further, SPINEX enables the generation of local explanations to gain insights into individual  
124 predictions. By considering the neighbors of a specific instance, the algorithm quantifies the  
125 importance of each feature and its interaction effects within the local context.

126 To facilitate the interpretation and analysis of SPINEX -based models, the proposed algorithm  
127 provides various visualization techniques, including feature importance plots, which show the  
128 relative contribution of each feature to the prediction, and interaction effect heatmaps, which  
129 visualize the interaction effects between different feature combinations. The algorithm also offers  
130 tools to analyze the change in predictions with the addition of neighbors, enabling the exploration  
131 of the model's behavior in response to nearby instances. Table 1 qualitatively compares SPINEX to  
132 other commonly used ML algorithms, which were also utilized in this work's experiments in a later  
133 section.

134 The following discussion further articulates the working mechanisms (i.e., functions and methods)  
135 of SPINEX for two versions of this algorithm: SPINEXRegressor and SPINEXClassifier.

136 Data Preprocessing: The SPINEX algorithm applies several preprocessing steps using the  
137 DataPreprocessor class. Given data matrix  $X \in \mathbb{R}^{n \times d}$  and corresponding label vector  $y \in \mathbb{R}^n$ ,  
138 preprocess the data to handle missing values, outliers, and perform feature selection. This includes  
139 handling missing data (removing or imputing), outlier detection and removal, and feature selection.  
140 The feature selection process uses a local search strategy, and the user can specify prioritized  
141 features for inclusion.

142 Distance Calculation: Distances between instances are calculated in the feature space. For a given  
143 distance metric  $d$ , calculate the distance  $D_{ij}$  between every pair of instances  $(x_i, x_j)$  as  $D_{ij} = d(x_i, x_j)$ .  
144 These distances are computed using the specified distance metric (Manhattan distance by default)  
145 and are then stored in a distance matrix. The user can specify the distance metric according to the  
146 problem at hand.

147 Weight Calculation: Weights are assigned to each of the training instances based on their distances  
148 to the test instance. The weights are computed using the Gaussian kernel function, where instances

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

149 closer to the test instance will have higher weights such that  $w_i$  for each instance  $x_i$  based on their  
150 distances to the test instance  $x_t$  as follows:  $w_i = \exp( - (D_{it}^2 / (2\sigma^2) ) )$ , where  $\sigma$  is the standard  
151 deviation of the distances. The standard deviation for the Gaussian kernel is computed as the mean  
152 of the distances.

153 Prediction: Predictions are made by considering the  $n$  nearest neighbors to a given test instance.  
154 The label  $\hat{y}_t$  for a test instance  $x_t$  as  $\hat{y}_t = \operatorname{argmax}_y \sum w_i * I(y_i = y)$ , where  $I$  is the indicator function  
155 that is 1 if  $y_i$  equals  $y$  and 0 otherwise, and the sum is over the  $n$  nearest neighbors of the test  
156 instance. The number of neighbors,  $n$ , is a user-defined parameter. The algorithm identifies these  
157 nearest neighbors based on the distance matrix. Once the neighbors are identified, they are  
158 combined based on the assigned weights to make a prediction.

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

159 Table 1 Qualitative comparison between SPINEX and commonly used ML algorithms  
160

Algorithm	Feature Importance	Interaction Effects	Interpretability	Model Complexity	Ensemble Learning	Handling Categorical Features	Handling Imbalanced Data	Speed	Inference	Model Size	Parallel Computing
SPINEX	Yes	Yes	Medium	Medium	Yes	No	No	Medium	Fast	Medium	Yes
Logistic Regression	Yes	No	High	Low	No	No	No	High	Fast	Small	Yes
Decision Tree	Yes	No	High	Varies	No	Yes	No	High	Fast	Varies	Yes
Random Forest	Yes	No	Medium	High	Yes	Yes	Yes	Medium	Fast	Large	Yes
Gradient Boosting	Yes	No	Medium	High	Yes	Yes	Yes	Low	Fast	Large	Yes
AdaBoost	Yes	No	Medium	High	Yes	Yes	Yes	Medium	Fast	Large	Yes
CatBoost	Yes	No	Medium	High	Yes	Yes	Yes	Medium	Fast	Large	Yes
XGBoost	Yes	No	Medium	High	Yes	Yes	Yes	High	Fast	Large	Yes
LightGBM	Yes	No	Medium	High	Yes	Yes	Yes	High	Fast	Large	Yes
Support Vector Classifier	No	No	Low	High	No	Yes	Yes	Low	Slow	Large	Yes
K-Nearest Neighbors	No	No	High	Low	No	No	No	Low	Slow	Small	Yes

- 161
- 162
- 163
- 164
- 165
- 166
- 167
- 168
- 169
- 170
- 171
- 172
- 173
- Feature Importance: Indicates whether the algorithm can provide feature importance or contribution scores.
  - Interaction Effects: Indicates whether the algorithm can capture and quantify interaction effects between features.
  - Interpretability: Assesses the ease of understanding and explaining the model's behavior and predictions.
  - Model Complexity: Indicates the complexity of the model. It assesses the level of complexity in terms of the number of parameters or rules used by the algorithm.
  - Ensemble Learning: Indicates whether the algorithm supports ensemble learning. Ensemble learning combines multiple models to improve performance.
  - Handling Categorical Features: Indicates whether the algorithm has built-in mechanisms to handle categorical features.
  - Handling Imbalanced Data: Indicates whether the algorithm has techniques to handle imbalanced datasets, where the number of instances in different classes is unequal.
  - Speed: Represents the speed of the algorithm for training and prediction tasks. It assesses the algorithm's efficiency in terms of computational time.
  - Inference: Indicates whether the algorithm supports efficient inference or prediction on new, unseen data after training.
  - Model Size: Assesses the memory footprint or storage requirements of the model.
  - Parallel Computing: Indicates whether the algorithm can leverage parallel computing capabilities to speed up training or prediction tasks.



Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

174 Feature contribution refers to the individual impact or importance of each feature on the prediction  
175 made by SPINEX. The algorithm calculates feature contributions using a neighbor-based approach.  
176 Once the neighbors are determined, the algorithm analyzes the difference in predictions between  
177 the original instance and its neighbors. This difference reflects the contribution of each feature to  
178 the prediction. Specifically, the algorithm compares the prediction made by the model when a  
179 particular feature is present in the instance with the prediction when that feature is absent.  
180 Mathematically, the contribution  $C_k$  of a feature  $f_k$  for an instance  $x_i$  as  $C_k = P(y|x_i) - P(y|x_i \wedge \neg f_k)$ ,  
181 where  $P(y|x_i)$  is the prediction probability for the instance  $x_i$  and  $P(y|x_i \wedge \neg f_k)$  is the prediction  
182 probability for the instance  $x_i$  with the  $k$ -th feature excluded. Calculate the interaction effect  $I_{kl}$   
183 between two features  $f_k$  and  $f_l$  as  $I_{kl} = C_k + C_l - C_{kl}$ , where  $C_{kl}$  is the change in prediction probability  
184 when both features are excluded.

185 The larger the difference, the more significant the contribution of the feature to the prediction. The  
186 feature contribution calculation takes into account the influence of neighboring instances, allowing  
187 the algorithm to capture the contextual importance of each feature. This approach provides a more  
188 nuanced understanding of feature importance, as it considers both the individual feature effects  
189 and their dependencies on nearby instances.

- 190 ○ In regression and classification algorithms, feature contributions are calculated by  
191 predicting the output with and without a given feature, then taking the difference. This  
192 indicates how much the prediction changes when a feature is removed, i.e., the  
193 "contribution" of the feature.
  - 194 ✓ In the regression algorithm, the method `compute_contributions()` is used to  
195 calculate feature contributions. It does this by predicting the output value with  
196 and without each feature and taking the difference.
  - 197 ✓ In the classification algorithm, the method `predict_contributions()` is used to  
198 calculate feature contributions. It does this by predicting class probabilities with  
199 and without each feature and taking the difference.
- 200 ○ Interaction effects refer to the combined impact of feature combinations on the  
201 prediction made by the SPINEX regression model. The algorithm analyzes the  
202 interactions between pairs or sets of features to identify synergistic or antagonistic  
203 effects that go beyond the individual contributions of the features. The algorithm  
204 examines the change in predictions when specific feature combinations are present or  
205 absent from calculating interaction effects. It compares the prediction made by the  
206 model with a particular feature combination to the prediction when that feature  
207 combination is removed. The difference in predictions reflects the interaction effect  
208 between the features in the combination. The algorithm considers all possible feature  
209 combinations, ranging from pairs to larger sets, to explore the full landscape of  
210 interactions. It quantifies the impact of each feature combination on the prediction,  
211 providing insights into the synergies or antagonisms between different features.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 212 ✓ In the regression algorithm, `compute_combination_impact()` is used to
- 213 calculate interaction effects. It does this by predicting the output value with all
- 214 features, with one feature removed and two removed, and then calculating the
- 215 interaction effect as described above.
- 216 ✓ In the classification algorithm, the method `predict_contributions()` is used to
- 217 calculate interaction effects. It does this by predicting class probabilities with
- 218 all features, with one feature removed and two features removed, and then
- 219 calculating the interaction effect described above.
- 220 ○ Overall, the calculation of feature contributions and interaction effects is very similar
- 221 in both the SPINEX classification and regression algorithms. The main difference is in
- 222 the predicted output (class probabilities vs. output value) and the context in which these
- 223 calculations are used (classification vs. regression).

224 Feature Importance and Impact Analysis: Feature importance is calculated as the mean absolute  
225 contribution of a feature across all instances. For example, the importance  $F_k$  of a feature  $f_k$  as  $F_k$   
226  $= (1/n) \sum |C_k|$ . Calculate the impact  $I_F$  of a combination of features  $F$  as  $I_F = P(y|x_i) - P(y|x_i^{-F})$ ,  
227 where  $P(y|x_i^{-F})$  is the prediction probability when all features in  $F$  are excluded. The impact of  
228 feature combinations is calculated by excluding a combination of features and computing the  
229 change in prediction probabilities. The results are sorted by impact, providing insight into which  
230 combinations of features are most important.

- 231 ■ Model Ensembling: The SPINEX model can be used as a base classifier in ensemble models.  
232 The user can specify an ensemble method, and the SPINEX classifier is then combined with  
233 other classifiers (like `DecisionTreeClassifier`) to make a final prediction. Three ensemble  
234 methods are used in the script: Stacking (where the predictions of the base classifiers are  
235 used as input to a final classifier), Bagging (where multiple instances of the SPINEX  
236 classifier are trained on random subsets of the training data), and Boosting (where multiple  
237 instances of the SPINEX classifier are trained sequentially, with each one focusing on the  
238 instances that the previous classifiers misclassified). For base classifiers  $h_1, \dots, h_p$ , in the  
239 case of Stacking, calculate the final prediction  $h(x)$  as  $h(x) = g(h_1(x), \dots, h_p(x))$ , where  $g$  is the  
240 final classifier. In the case of Bagging, calculate  $h(x)$  as  $h(x) = \text{majority}(h_1(x), \dots, h_p(x))$ . In  
241 the case of Boosting, calculate  $h(x)$  as  $h(x) = \text{weighted\_majority}(h_1(x), \dots, h_p(x))$ .

### 242 2.3 SPINEX for regression

#### 243 Inputs:

244 `X_train`: Training feature matrix  
245 `y_train`: Training target vector  
246 `X_test`: Test feature matrix  
247 `distance_metric`: Distance metric for calculating pairwise distances  
248 `num_neighbors`: Number of nearest neighbors to consider  
249 `kernel_width`: Width of the Gaussian kernel  
250



This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

```
251 Procedure SPINEX:
252   Preprocess(X_train, X_test)
253   distances = CalculateDistances(X_train, X_test)
254   weights = CalculateWeights(distances)
255   predictions = Predict(X_train, y_train, X_test, weights)
256   Return predictions
257
258 Procedure Preprocess(X_train, X_test):
259   # Handle missing data and outliers in X_train and X_test
260
261 Procedure CalculateDistances(X_train, X_test):
262   distances = empty matrix of size (number of test instances) x (number of training instances)
263   For each test_instance in X_test:
264     For each train_instance in X_train:
265       distances[test_instance][train_instance] = calculateDistance(test_instance, train_instance, distance_metric)
266   Return distances
267
268 Procedure CalculateWeights(distances):
269   weights = empty matrix of size (number of test instances) x (number of training instances)
270   For each test_instance in distances:
271     sorted_distances = sort(distances[test_instance]) # Sort distances in ascending order
272     kernel_bandwidth = kernel_width * mean(sorted_distances) # Compute kernel bandwidth
273     For i = 0 to num_neighbors - 1:
274       weights[test_instance][i] = calculateWeight(sorted_distances[i], kernel_bandwidth)
275   Return weights
276
277 Procedure Predict(X_train, y_train, X_test, weights):
278   predictions = empty vector of size (number of test instances)
279   For each test_instance in X_test:
280     nearest_neighbors = GetNearestNeighbors(weights[test_instance], num_neighbors)
281     prediction = CalculatePrediction(nearest_neighbors, y_train)
282     predictions[test_instance] = prediction
283   Return predictions
284
285 Procedure GetNearestNeighbors(weights, num_neighbors):
286   sorted_indices = indices of weights sorted in descending order
287   nearest_neighbors = first num_neighbors indices from sorted_indices
288   Return nearest_neighbors
289
290 Procedure CalculatePrediction(nearest_neighbors, y_train):
291   prediction = average of y_train values corresponding to nearest_neighbors
292   Return prediction
```

## 293 Explanation of main functions and methods:

- 294 ■ Fitting and Predicting with SPINEXRegressor

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 295 - fit(X, y): This method trains the SPINEX regression model using training data X and target  
296 values y.
- 297 - predict(X): This method predicts the target values for a given set of input samples X, using  
298 the trained SPINEX regression model.
- 299 - predict\_contributions(X): This method predicts the contributions of each feature to the  
300 target values for a given set of input samples X.
- 301 ■ Analysis & Visualization of Feature Importance and Interaction Effects
- 302 - get\_feature\_importance(X): This method calculates the feature importances for each  
303 feature in X.
- 304 - get\_global\_interaction\_effects(X): This method calculates the average interaction effects  
305 for each feature in the dataset.
- 306 - feature\_combination\_impact\_analysis(X): This method analyzes the impact of different  
307 combinations of features on the model's predictions.
- 308 - normalize\_importances(importances): A utility function for normalizing feature  
309 importances.
- 310 - visualize\_feature\_importances(local\_importances, global\_importances,  
311 feature\_names): This method generates a bar plot to compare local and global  
312 feature importances.
- 313 - visualize\_interaction\_effects(interaction\_effects\_df): This method generates a bar plot to  
314 show the interaction effects between different features.
- 315 - plot\_average\_interaction\_network(avg\_interaction\_effects, feature\_names=None):  
316 This function creates a network graph to visualize the interactions between different  
317 features.
- 318 - plot\_contribution\_heatmaps(contributions, interaction\_effects, feature\_names=None):  
319 This function creates two heatmaps - one for individual feature contributions and  
320 one for pairwise interactions.
- 321 ■ Local Explanations & Visualizations
- 322 - get\_local\_explanation(X, instance\_to\_explain): This method calculates the local feature  
323 importances for a specific instance.
- 324 - get\_local\_interaction\_effects(X, instance\_to\_explain): This method calculates the local  
325 interaction effects for a specific instance.
- 326 - plot\_prediction\_change(X, y, instance\_to\_explain): This function visualizes how the  
327 prediction changes as each neighbor is added.
- 328 - visualize\_neighbor\_counts(neighbor\_counts): This function visualizes the counts of each  
329 neighbor in a bar plot.
- 330 ■ Influence of Feature Combinations & Local Changes
- 331 - plot\_feature\_pair\_influence(X, instance\_to\_explain, feature\_pair, grid\_size=20): This  
332 method generates a 3D plot to show how changing the values of a pair of features  
333 influences the prediction.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 334 - plot\_feature\_triplet\_influence(X, instance\_to\_explain, feature\_triplet, feature\_names,  
335 grid\_size=20): This method generates a 3D scatter plot to show how changing the  
336 values of a triplet of features influences the prediction.
- 337 - explore\_local\_changes(X, instance\_to\_explain, feature\_to\_explore, grid\_size=100,  
338 feature\_range=None): This method generates a series of predictions by varying the  
339 value of a specific feature and keeping all other features constant.
- 340 - plot\_local\_changes(X, instance\_to\_explain, feature\_to\_explain, grid\_size=100,  
341 feature\_range=None, feature\_name=None): This function visualizes how the  
342 prediction changes as the value of a specific feature changes.
- 343 - explore\_all\_local\_changes(X, instance\_to\_explain, grid\_size=100, feature\_range=None):  
344 This function generates a series of predictions by varying the values of all features  
345 one by one and keeping all other features constant.
- 346 - explore\_local\_changes\_for\_pair(X, instance\_to\_explain, feature\_pair, grid\_size=100,  
347 feature\_range=None): This method generates a series of predictions by varying the  
348 values of a pair of features and keeping all other features constant.
- 349 - explore\_local\_changes\_for\_triplet(X, instance\_to\_explain, feature\_triplet, grid\_size=10,  
350 feature\_range=None): This method generates a series of predictions by varying the  
351 values of a triplet of features and keeping all other features constant.

352 The following are the hyperparameters for the regression version of SPINEX (many of which are  
353 similar to those for the classification version):

- 354 ■ n\_neighbors: This parameter controls the number of neighbors to use for neighbor queries.  
355 The choice of n\_neighbors affects the predictions made by the model: a smaller number  
356 makes the model more sensitive to local variations in the data, while a larger number makes  
357 the predictions more stable at the expense of potentially ignoring smaller patterns.
- 358 ■ distance\_threshold: This parameter is used in the calculation of instance weights. Weights  
359 are calculated as the reciprocal of the sum of the distance to each neighbor and the decayed  
360 distance threshold. The distance\_threshold parameter thus controls how much influence  
361 more distant neighbors have on the prediction of a given instance.
- 362 ■ distance\_threshold\_decay: This parameter controls the decay rate of the distance  
363 threshold. A lower decay rate means that the influence of more distant neighbors decays  
364 more quickly.
- 365 ■ ensemble\_method: This parameter allows the user to specify an ensemble method to use  
366 in combination with SPINEX. Options include bagging, boosting, and stacking. Ensemble  
367 methods combine the predictions of multiple models to improve predictive performance.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 368       ▪ `n_features_to_select`: This parameter controls the number of features to select for training  
369       the model. If `auto_select_features` is set to `True`, the model will automatically select the  
370       features it deems most important.
- 371       ▪ `auto_select_features`: If this parameter is set to `True`, the model will automatically select  
372       a subset of features for training. The number of features to select is controlled by the  
373       `n_features_to_select` parameter.
- 374       ▪ `use_local_search`: If this parameter is set to `True`, the model will perform a local search to  
375       select the best features for training. This can potentially improve the model's performance  
376       but may also increase training time.
- 377       ▪ `prioritized_features`: This parameter allows the user to specify a list of features that should  
378       be prioritized in the feature selection process.
- 379       ▪ `missing_data_method`: This parameter allows the user to specify the method for handling  
380       missing data. Options include `mean_imputation`, which replaces missing values with the  
381       mean of the existing values, and `deletion`, which removes instances with missing values.
- 382       ▪ `outlier_handling_method`: This parameter allows the user to specify the method for  
383       handling outliers. Options include `z_score_outlier_handling`, which removes instances that  
384       have a Z-score greater than 3, and `iqr_outlier_handling`, which removes instances that fall  
385       outside a certain range defined by the interquartile range (IQR).
- 386       ▪ `exclude_method`: This parameter allows the user to specify a method for excluding certain  
387       instances from the training set. This could be used, for example, to exclude instances  
388       considered outliers based on some criterion.

## 389 *2.4 SPINEX for classification*

### 390 **Inputs:**

391 `X_train`: Training feature matrix  
392 `y_train`: Training target vector (class labels)  
393 `X_test`: Test feature matrix  
394 `distance_metric`: Distance metric for calculating pairwise distances  
395 `num_neighbors`: Number of nearest neighbors to consider  
396 `kernel_width`: Width of the Gaussian kernel  
397

### 398 **Procedure SPINEX:**

399 `Preprocess(X_train, X_test)`  
400 `distances = CalculateDistances(X_train, X_test)`  
401 `weights = CalculateWeights(distances)`  
402 `predictions = Predict(X_train, y_train, X_test, weights)`  
403 `Return predictions`  
404

### 405 **Procedure Preprocess(X\_train, X\_test):**

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

```
406 # Handle missing data and outliers in X_train and X_test
407
408 Procedure CalculateDistances(X_train, X_test):
409     distances = empty matrix of size (number of test instances) x (number of training instances)
410     For each test_instance in X_test:
411         For each train_instance in X_train:
412             distances[test_instance][train_instance] = calculateDistance(test_instance, train_instance, distance_metric)
413     Return distances
414
415 Procedure CalculateWeights(distances):
416     weights = empty matrix of size (number of test instances) x (number of training instances)
417     For each test_instance in distances:
418         sorted_distances = sort(distances[test_instance]) # Sort distances in ascending order
419         kernel_bandwidth = kernel_width * mean(sorted_distances) # Compute kernel bandwidth
420         For i = 0 to num_neighbors - 1:
421             weights[test_instance][i] = calculateWeight(sorted_distances[i], kernel_bandwidth)
422     Return weights
423
424 Procedure Predict(X_train, y_train, X_test, weights):
425     predictions = empty vector of size (number of test instances)
426     For each test_instance in X_test:
427         nearest_neighbors = FindNearestNeighbors(weights[test_instance], num_neighbors)
428         prediction = CalculatePrediction(nearest_neighbors, y_train)
429         predictions[test_instance] = prediction
430     Return predictions
431
432 Procedure FindNearestNeighbors(weights, num_neighbors):
433     sorted_indices = indices of weights sorted in descending order
434     nearest_neighbors = first num_neighbors indices from sorted_indices
435     Return nearest_neighbors
436
437 Procedure CalculatePrediction(nearest_neighbors, y_train):
438     class_counts = empty dictionary
439     For each neighbor in nearest_neighbors:
440         class_label = y_train[neighbor]
441         If class_label is not in class_counts:
442             class_counts[class_label] = 1
443         Else:
444             class_counts[class_label] += 1
445     prediction = class label with the highest count in class_counts
446     Return prediction
```

#### 447 Explanation of main functions and methods:

- 448 ■ Fitting and Predicting with SPINEXClassifier
- 449 - fit(X, y): This method trains the SPINEX classification model using training data X and
- 450 target values y.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 451           - predict(X): This method predicts the target values for a given set of input samples X, using
- 452                   the trained SPINEX classification model.
- 453       ▪ Assessing Model Accuracy
- 454           - score(X, y): This method calculates the mean accuracy on the given test data and labels.
- 455       ▪ Analysis & Visualization of Feature Importance and Interaction Effects
- 456           - get\_feature\_contributions(X): This method calculates the contributions of each feature
- 457                   for a given instance.
- 458           - plot\_feature\_contributions(feature\_contributions): This method generates a bar plot to
- 459                   visualize the contributions of each feature.
- 460           - get\_feature\_interactions(X): This method calculates the interactions between every pair
- 461                   of features for a given instance.
- 462           - plot\_feature\_interactions(feature\_interactions): This method generates a heatmap to
- 463                   visualize the interaction effects between different features.
- 464           - plot\_prediction\_change(X, y): This function visualizes how the prediction changes as
- 465                   each neighbor is added.
- 466       ▪ Local Explanations & Visualizations
- 467           - get\_local\_explanation(X, instance\_to\_explain): This method calculates the local feature
- 468                   importances and neighbor counts for a specific instance.
- 469           - get\_local\_interaction\_effects(X, instance\_to\_explain): This method calculates the local
- 470                   interaction effects for a specific instance.
- 471           - visualize\_neighbor\_counts(neighbor\_counts): This function visualizes the counts of each
- 472                   neighbor in a bar plot.
- 473       ▪ Influence of Feature Combinations & Local Changes
- 474           - plot\_all\_feature\_contributions(model, X): This function generates a scatter plot to
- 475                   visualize the contributions of all features.
- 476           - get\_global\_interaction\_effects(X, y) and feature\_combination\_impact\_analysis(X, y):
- 477                   These functions calculate the average interaction effects and the impact of feature
- 478                   combinations for all instances in the dataset.
- 479           - get\_feature\_importance(X, instance\_to\_explain): This method obtains the feature
- 480                   importances and interaction effects for selected instances.
- 481           - plot\_feature\_pair\_influence(X, instance\_to\_explain, feature\_pair): This method
- 482                   generates a scatter plot to show how changing the values of a pair of features
- 483                   influences the prediction.
- 484           - plot\_feature\_triplet\_influence(X, instance\_to\_explain, feature\_triplet): This method
- 485                   generates a scatter plot to show how changing the values of a triplet of features
- 486                   influences the prediction.

### 487 **3.0 Description of benchmarking experiments, algorithms, and datasets**

488 This section describes the experimental examination used to benchmark SPINEX. For a start, SPINEX  
489 was examined against ten other commonly used ML algorithms, namely, Logistic Regression,



Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

490 Decision Tree, Random Forest, Gradient Boosting, AdaBoost, CatBoost, XGBoost, LightGBM,  
491 Support Vector Classifier, and K-Nearest Neighbors. All of these models were used in their default  
492 settings<sup>1</sup>.

493 A series of synthetic and real regression and classification datasets were used. Many of these  
494 datasets were also recently benchmarked via several ML algorithms [34,35]. Each dataset is  
495 checked per the recommendations of recent researchers aimed at measuring data health. Three  
496 criteria were selected, and all of these criteria were satisfied,

- 497 • Van Smeden et al. [36] require a minimum set of 10 observations per feature.
- 498 • Riley et al. [37] suggest a minimum of 23 cases per feature.
- 499 • Frank and Todeschini [38] recommend maintaining a minimum ratio of 3 and 5 between  
500 the number of observations and features.

501 A 5-fold cross-validation technique is applied in regression experiments, and in classification  
502 experiments, a stratified 10-fold cross-validation technique is applied [39–41]. The performance  
503 of the ML models created is evaluated via a number of regression and classification metrics, as  
504 listed in Table 2 [42]. For regression problems, the metrics included the mean absolute error  
505 (MAE) and coefficient of determination ( $R^2$ ). In general, lower values of MAE and values close  
506 to positive unity for  $R^2$  are favorable. In addition, the classification metrics include accuracy,  
507 logloss error, and the area under the receiver operating characteristic (ROC) curve (AUC).  
508 Naturally, higher values of accuracy and AUC and lower values of the logloss metrics are  
509 favorable. Finally, a newly-derived functional metrics is also used estimated energy.

510 Table 2 List of common performance metrics.

Metric	Formula
Regression	
Mean Absolute Error (MAE)	$MAE = \frac{\sum_{i=1}^n  P_i - A_i }{n}$
Coefficient of Determination ( $R^2$ )	$R^2 = 1 - \frac{\sum_{i=1}^n (P_i - A_i)^2}{\sum_{i=1}^n (A_i - A_{mean})^2}$ <i>A: actual measurements, P: predictions, n: number of data points.</i>
Classification	
Accuracy	$ACC = \frac{TP + TN}{P + N}$ <i>P: predictions, N: number of real negatives, TP: number of true positives, TN: number of true negatives.</i>

<sup>1</sup> The default setting for SPINEX include:

- SPINEXRegressor = SPINEX(n\_neighbors=5, distance\_threshold=0.05, distance\_threshold\_decay=0.05, ensemble\_method=None, n\_features\_to\_select=None, auto\_select\_features=False, use\_local\_search=False, prioritized\_features=None, missing\_data\_method='mean\_imputation', outlier\_handling\_method='z\_score\_outlier\_handling', exclude\_method='zero')
- SPINEXClassifier = SPINEX(n\_neighbors=5, distance\_threshold=0.05, distance\_threshold\_decay=0.95, ensemble\_method=None, metric='euclidean')

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

Logloss error	$LLE = - \sum_{c=1}^M A_i \log P$ <p><i>M</i>: number of classes, <i>c</i>: class label, <i>y</i>: binary indicator (0 or 1) if <i>c</i> is the correct classification for a given observation.</p>
Area under the ROC curve	$AUC = \sum_{i=1}^{N-1} \frac{1}{2} (FP_{i+1} - FP_i) (TP_{i+1} - TP_i)$ <p><i>FP</i>: number of false positives, <i>FN</i>: number of false negatives.</p>
Functional metric (from [34])	
Estimated Energy	$MS \times (TT + PT)$ <p><i>MS</i>: model size, <i>TT</i>: training time, and <i>PT</i>: prediction time. Smaller values are favorable with a hypothetical minimum value = 1.0 MB × 10 sec = 10 MB.sec.</p>

### 511 3.1 Synthetic datasets

512 A collection of functions that generate synthetic data sets, each simulating different types of  
 513 relationships between features and the target, were used in our experiments. These datasets are  
 514 useful for testing and evaluating ML models' performance (see Table 3).

#### 515 3.1.1 Regression experiments

516 The `make_regression` function from the `sklearn.datasets` module was used herein. More  
 517 specifically, four primary functions are used. These include:

- 518 ■ `generate_regression_data` function generates random regression data. The underlying  
 519 equation for this function can be represented as:

$$520 y = X_1 w_1 + X_2 w_2 + \dots + X_n w_n + b + \epsilon$$

521 Here, 'X' represents the input features, *w* signifies the weights, *b* is the bias, and  $\epsilon$  is a random  
 522 noise term.

- 523 ■ `generate_synthetic_data` function generates a dataset where the relationship between the  
 524 features and the target is a combination of a quadratic function, a sinusoidal function, and a  
 525 simple multiplication. This is a somewhat complex relationship, which would provide a  
 526 challenge to many types of machine learning models. The equation for this function is:

$$527 y = X_0 + X_1^1 + X_2^2 + \dots + X_n^n + \epsilon + \text{outlier\_noise}$$

528 In this equation,  $\epsilon$  represents random noise, while `outlier_noise` is additional noise added to  
 529 randomly chosen samples.

- 530 ■ `generate_cubic_data` function creates a dataset where the relationship between the features  
 531 and the target is a combination of a cubic function, a squared function, and a simple  
 532 multiplication. This type of dataset is useful for testing how well a model can handle cubic  
 533 relationships. Here, the equation is:

$$534 y = X_0 + X_1^3 + X_2^4 + \dots + X_n^{n+2} + \epsilon + \text{outlier\_noise}$$

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

535     ▪ generate\_exponential\_data function creates a dataset with an exponential relationship  
536     between one of the features and the target, a decaying exponential for another feature, and  
537     simple linear relationships for the remaining features. The equation for this function is:

538     
$$y = \exp(X_0) + X_1^1 + X_2^2 + \dots + X_n^n + \epsilon + \text{outlier\_noise}$$

539     ▪ generate\_step\_data function creates synthetic data where the relationship between the first  
540     feature and the output is a step function, and subsequent features contribute polynomially to  
541     the output. Gaussian noise is added to the output. The underlying equation can be represented  
542     as:

543     
$$y = u(X_0 - 0.5) + X_1^1 + X_2^2 + \dots + X_n^n + \epsilon$$

544     In this equation,  $u$  is the unit step function.

545     ▪ generate\_complex\_interaction\_data function generates synthetic data with complex  
546     interactions between the features, including polynomial, sinusoidal, and logarithmic  
547     interactions. Here, the equation is:

548     
$$y = X_0^2 + \sin(X_1) \times \log(X_2^2 + 1) + \epsilon$$

549     ▪ generate\_polynomial\_data function creates synthetic data where the relationship between the  
550     features and the output is defined by high degree polynomials. The underlying equation can be  
551     represented as:

552     
$$y = X_0^3 + X_1^4 - X_2^5 + \epsilon$$

553     ▪ generate\_exp\_log\_data function generates an interaction between an exponential function of  
554     the first feature and the natural logarithm of one plus the second feature. The equation for this  
555     function is:

556     
$$y = \exp(X_0) \times \log_{10}(X_1) + \epsilon$$

557     ▪ generate\_sin\_exp\_data function creates synthetic data where the output is an interaction  
558     between a sinusoidal function of the first feature and an exponential function of the second  
559     feature. Here, the equation is:

560     
$$y = \sin(\pi X_0) \times \exp(X_1) + \epsilon$$

561     ▪ generate\_tan\_data function generates a tangent of the first feature, with subsequent features  
562     contributing polynomially to the output. The underlying equation can be represented as:

563     
$$y = \tan(X_0) + X_1^1 + X_2^2 + \dots + X_n^n + \epsilon$$

564     These functions accepts parameters such as:

565     ○ n\_samples for the number of samples

566     ○ n\_features for the number of features

567     ○ n\_informative for the number of informative features, i.e., the features that are useful in  
568     predicting the target variable. The remaining features ( $n\_features - n\_informative$ ) are  
569     generated as random noise.

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 570 ○ noise for the standard deviation of the Gaussian noise applied to the output (dependent  
571 variable).
- 572 ○ n\_targets for the number of regression targets, i.e., the number of dependent variables. By  
573 default, this is set to 1, meaning that the generated dataset will have a single target variable.
- 574 ○ n\_outliers for the number of outliers.
- 575 ○ bias for the bias term in the underlying linear model.
- 576 ○ shuffle to assign whether or not to shuffle the samples and the features.
- 577 ○ effective\_rank for the approximate number of singular vectors required to explain most of  
578 the input data by a linear, low rank model. If None, all features are informative. This  
579 parameter can be used to introduce collinearity in the data.
- 580 ○ tail\_strength for the relative importance of the fat noisy tail of the singular values profile  
581 if effective\_rank is not None.
- 582 ○ seed to assign seed for the random number generator, to ensure the reproducibility of the  
583 results.

584 In all functions, combination of parameters was assigned to create 18 datasets (see Table 3). The  
585 amount of noise can be increased or decreased by adjusting the noise\_scale parameter. In addition,  
586 each function adds a specified number of outliers to the target variable. The outliers are randomly  
587 selected from the samples and have a larger amount of normally distributed random noise added  
588 to them.

589 Finally, a dictionary named datasets is created to store the generated synthetic datasets. Each  
590 dataset has a descriptive name and is generated by one of the previously described functions, with  
591 specific parameters. This dictionary allows easy access to each dataset by its name, which is useful  
592 when looping through them later to fit and evaluate different ML models.

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

593 Table 3 Datasets used in the regression analysis on synthetic data.

Dataset	Function	n_samples	n_features	n_informative	n_outliers	noise	bias	shuffle	effective_rank	tail_strength
Dataset 1	generate_regression_data	50	5	5	-	0.0	0.0	-	-	-
Dataset 2	generate_regression_data	5000	4	4	-	0.1	0.0	-	-	-
Dataset 3	generate_regression_data	1000	6	5	-	0.0	10	-	-	-
Dataset 4	generate_regression_data	7000	2	2	-	0.0	0.0	False	-	-
Dataset 5	generate_regression_data	750	8	6	-	0.0	0.0	-	5	-
Dataset 6	generate_regression_data	800	4	4	-	0.0	0.0	-	-	0.1
Dataset 7	generate_regression_data	1000	5	3	-	0.0	10	-	-	-
Dataset 8	generate_regression_data	2500	3	2	-	0.0	0.0	False	-	-
Dataset 9	generate_regression_data	1000	4	4	-	0.9	0.0	-	10	-
Dataset 10	generate_step_data	2000	7	-	-	0.0	-	-	-	-
Dataset 11	generate_cubic_data	1000	10	-	20	0.5	-	-	-	-
Dataset 12	generate_synthetic_data	2000	6	-	200	0.8	-	-	-	-
Dataset 13	generate_exponential_data	2000	5	-	40	0.8	-	-	-	-
Dataset 14	generate_tan_data	750	8	-	-	0.1	-	-	-	-
Dataset 15	generate_complex_interaction_data	500	7	-	-	0.0	-	-	-	-
Dataset 16	generate_polynomial_data	2000	5	-	-	0.1	-	-	-	-
Dataset 17	generate_exp_log_data	1000	10	-	-	0.5	-	-	-	-
Dataset 18	generate_sin_exp_data	3000	5	-	-	0.0	-	-	-	-

594

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

595 A systematic approach to rank models based on the selected multiple metrics is followed. In this  
596 approach, we calculate the average scores for each model across different metrics, assign ranks to  
597 the models for each metric, calculate the sum-based rank across all metrics, and display the ranked  
598 models. This enables a comparative analysis of models based on their performance across various  
599 metrics. The outcome of this analysis is shown in Table 4 as well as Fig. 1.

600 It is quite clear that SPINEX and most of its derivatives rank well in terms of accuracy and the  
601 bottom half of the total time for training and prediction. The rankings seem to fall in terms of total  
602 time (which also affect the ranking for energy). This is due to the algorithm's design to check for  
603 feature pairs and interactions.

604 Table 4 Ranking results of regression experiment on synthetic data

Model	MAE	R <sup>2</sup>	Rank	Total Time	Estimated Energy	Rank
StackingSPINEX	1	1	1	15	16	16
CatBoostRegressor	6	2	2	16	14	15
BayesianRidge	3	6	3	4	4	4
HuberRegressor	2	8	4	7	7	7
GradientBoostingRegressor	7	3	4	12	12	12
Ridge	4	7	5	2	1	1
XGBRegressor	8	5	7	9	11	10
RandomForestRegressor	9	4	6	14	15	14
LGBMRegressor	12	9	7	8	8	8
BaggingSPINEX	11	10	8	17	17	17
Lasso	5	17	9	1	2	1
SPINEX	13	11	10	13	13	13
BoostingSPINEX	10	15	11	18	18	18
KNeighborsRegressor	14	12	12	5	5	5
AdaBoostRegressor	15	14	13	10	9	9
SVR	18	13	14	11	10	11
DecisionTreeRegressor	17	16	15	6	6	6
ElasticNet	16	18	16	3	3	3

605

### 606 3.1.2 Classification experiments

607 Similar to the regression counterpart, a number of synthetic datasets for binary classification tasks  
608 were generated using the `make_classification` function from the `sklearn.datasets` module. The  
609 function `generate_synthetic_data` is defined to create synthetic datasets and accepts the following  
610 parameters:

- 611 ○ `n_samples` is the total number of data points in the dataset.



Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 612 ○ `n_features` is the total number of features in the dataset.
- 613 ○ `n_informative` is the number of informative features (i.e., useful for classifying the
- 614 samples).
- 615 ○ `n_redundant` is the number of redundant features (i.e., generated as random linear
- 616 combinations of the informative features).
- 617 ○ `weights` is the proportions of samples assigned to each class.
- 618 ○ `flip_y` is the fraction of samples whose classes are randomly exchanged.
- 619 ○ `class_sep` is the factor multiplying the hypercube size, wherein larger values spread out the
- 620 classes tend to make the classification task easier.

621 The function `make_classification` generates a random `n`-class classification problem. It returns `X`

622 and `y`, where `X` is a 2D array of shape `n_samples, n_features` representing the generated samples,

623 and `y` is a 1D array of shape `n_samples` representing the integer labels for class membership of

624 each sample. Overall, 18 datasets were synthetically generated and tested. These were labeled

625 under Series A (see Table 5), and Series B (see Table 5), with varying complexities.

626 A similar systematic approach to rank the ML models based on the selected multiple metrics is

627 followed here as that in the regression analysis. The outcome of this analysis is shown in Tables 6

628 and 7 as well as Figs. 2 and 3. Naturally, most models' predictions are slightly degraded when

629 evaluated on the datasets belonging to Series B (given their complexity).

630 The outcome of the analysis also shows that SPINEX and its derivatives perform much better in this

631 classification task than in the regression. Overall, and despite its relatively poor ranking under time

632 and energy consumption, the default version of SPINEX consistently ranks in the top 7 in the overall

633 ranking. It is clear that the SPINEX model can outperform some of the more common and traditional

634 algorithms, even in scenarios of imbalanced data and relatively large datasets.

635 Table 6 Ranking results of classification experiment on synthetic data (Series A)

Models	Accuracy	LLE	AUC	Rank	Estimated Energy	Total Time	Rank
SVC	1	1	1	1	10	15	13
SPINEXClassifier(default)	2	11	2	2	13	9	12
StackingSPINEX	7	2	7	2	14	12	14
BaggingSPINEX	4	7	5	4	15	14	15
KNeighborsClassifier	3	12	3	4	2	2	2
BoostingSPINEX	6	5	6	6	11	5	7
SPINEX	5	6	4	7	12	7	9
LGBMClassifier	8	4	10	8	4	3	3
XGBClassifier	9	3	9	9	6	6	5
RandomForestClassifier	10	8	8	10	8	10	8
GradientBoostingClassifier	11	9	11	11	9	11	10

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

AdaBoostClassifier	12	13	12	12	5	8	6
DecisionTreeClassifier	13	14	14	12	3	4	3
LogisticRegression	14	10	13	14	1	1	1
CatBoostClassifier	-	-	-	-	7	13	10

\*SPINEX = (n\_neighbors=20, distance\_threshold=0.05, distance\_threshold\_decay=0.95, ensemble\_method=None, metric='manhattan')

636

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

637 Table 5 Datasets used in the classification analysis on synthetic data

Dataset	Series A				Series B						
	n_samples	n_features	n_informative	n_redundant	n_samples	n_features	n_informative	n_redundant	flip_y	class_sep	weights
Dataset 1	50	3	2	0	50	3	2	0	0.01	1.0	0.9/0.1
Dataset 2	100	10	6	2	100	10	6	2	0.02	0.5	0.8/0.2
Dataset 3	1000	80	20	40	1000	80	20	40	0.03	0.8	0.7/0.3
Dataset 4	500	20	20	0	500	20	20	0	0.04	0.2	0.6/0.4
Dataset 5	5000	40	15	10	5000	40	15	10	0.05	0.3	0.5/0.5
Dataset 6	10000	10	5	5	10000	10	5	5	0.06	0.4	0.6/0.4
Dataset 7	500	20	20	0	1500	100	40	0	0.07	0.5	0.7/0.3
Dataset 8	3000	55	20	20	3000	55	20	20	0.08	0.6	0.8/0.2
Dataset 9	50000	5	3	0	50000	5	3	0	0.09	0.7	0.6/0.4

638

639 Table 7 Ranking results of classification experiment on synthetic data (Series B)

Models	Accuracy	LLE	AUC	Rank	Estimated Energy	Total Time	Rank
SPINEXClassifier(default)	3	1	1	1	13	9	12
StackingSPINEX	2	3	3	2	14	12	14
SPINEX	7	2	2	3	12	7	9
KNeighborsClassifier	4	5	5	4	2	2	2
BaggingSPINEX	6	4	4	4	15	15	15
LogisticRegression	14	6	6	6	1	1	1
RandomForestClassifier	11	8	8	7	9	11	10
DecisionTreeClassifier	15	7	7	8	3	4	3
SVC	1	14	14	8	8	10	8
CatBoostClassifier	8	11	11	10	7	14	11
GradientBoostingClassifier	12	9	9	10	10	13	13
AdaBoostClassifier	13	10	10	12	5	8	6
XGBClassifier	9	12	12	12	6	6	5
BoostingSPINEX	5	15	15	14	11	5	7
LGBMClassifier	10	13	13	15	4	3	3

640

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

641 **3.2 Real datasets**

642 Now, we repeat the analysis using a series of real datasets. These datasets are described in the  
643 following sections.

644 **3.2.1 Regression experiments**

645 Here, 12 real datasets are used to benchmark SPINEX and the other ML models. Table 8 lists details  
646 on each dataset, along with their respective references. As one can see, the selected datasets  
647 comprise a collection of samples and features and cover various problem domains. Further  
648 information can be found in each dataset’s reference.

649 **Table 8 Datasets used in the regression analysis on synthetic data**

Dataset	n_samples	n_features	Ref.
University Admission	401	8	[43]
Fire Resistance of RC columns	311	13	[44]
Shear Strength of beams	168	7	[35]
Concrete Strength	1031	9	[45]
Deformation of Beams under Fire	1187	7	[46]
Strength of Steel Tubes	1260	6	[47]
Energy Efficiency of Buildings	767	10	[48]
Body Fat Index	252	15	[49]
Forest Fire Area	517	13	[50]
Abalone Age	2000	10	[51]
Synchronous Motor	557	5	[52]
Walmart Retail	3000	6	[53]

650

651 A comparative analysis of models based on their performance across various metrics is presented  
652 in Table 9 as well as Fig. 4. The top performing SPINEX derivative ranks 3<sup>rd</sup> and 4<sup>th</sup> in terms of  
653 accuracy. Other SPINEX derivates also faired well in terms of accuracy metrics (and outperforming  
654 some of the common algorithms such as LGBMRegressor, RandomForestRegressor, and  
655 GradientBoostingRegressor, but continue to rank at the bottom half due to their large time and  
656 energy used.

657 **Table 9 Ranking results of regression experiment on synthetic data**

Model	MAE	R <sup>2</sup>	Rank	Total Time	Estimated Energy	Rank
CatBoostRegressor	1	1	1	3	3	2
BaggingSPINEX	8	2	2	14	15	14
SPINEX	4	7	3	13	14	12
LGBMRegressor	9	4	4	4	4	3
RandomForestRegressor	3	10	5	12	12	11
StackingSPINEX	11	3	6	16	16	15

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

GradientBoostingRegressor	5	12	7	15	13	13
HuberRegressor	12	6	8	6	6	5
AdaBoostRegressor	13	5	9	8	10	8
XGBRegressor	2	16	10	10	11	10
KNeighborsRegressor	10	11	11	5	5	4
BayesianRidge	14	8	12	1	2	1
Ridge	15	9	13	2	1	1
DecisionTreeRegressor	7	17	14	9	8	7
BoostingSPINEX	6	18	15	18	18	17
Lasso	16	14	16	7	7	6
SVR	18	13	17	17	17	16
ElasticNet	17	15	18	11	9	9

658

### 659 3.2.2 Classification experiments

660 Here, 11 real datasets are used to examine SPINEX and the other ML models as a means for a second  
 661 means of validation. Table 10 lists details for each used dataset regarding the number of samples  
 662 and features. Further information can be found in each dataset's reference.

663 **Table 10 Datasets used in the classification analysis on real data**

Dataset	n_samples	n_features	Ref.
Fire-induced Spalling	1062	16	[54]
Pima Indians Diabetes	768	8	[55]
Bridge Failures	299	7	[56]
Concrete Condition in Situ	9683	8	[57]
Breast Cancer Wisconsin (Original)	569	30	[58,59]
Rice (Commeo and Osmancik)	3810	7	[60]
Bank Note Authentication	1372	4	[61]
Water Potability	2011	9	[62]
Machine Predictive Maintenance	10000	5	[63]
Depression Prediction	1409	20	[64]
Cars Purchase Decision	1000	3	[65]

664

665 A look into Table 11 shows that GradientBoostingClassifier and CatBoostClassifier seem to rank  
 666 constantly among the top two models in terms of accuracy. On the other hand, three versions of  
 667 SPINEX (namely, StackingSPINEX and SPINEX) land at 6<sup>th</sup> and 7<sup>th</sup> in the overall ranking for accuracy.  
 668 Figure 5 shows that despite the low ranking of SPINEX, this algorithm scored comparable  
 669 performance to other traditional ML models, such as KNeighborsClassifier, DecisionTreeClassifier,  
 670 LogisticRegression, and SVC.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

671 Table 11 Ranking results of classification experiment on real data

Models	Accuracy	LLE	AUC	Rank	Estimated Energy	Total Time	Rank
RandomForestClassifier	3	5	3	1	9	13	9
GradientBoostingClassifier	1	2	2	2	8	12	7
CatBoostClassifier	2	1	1	3	10	15	10
XGBClassifier	5	4	5	3	6	8	5
LGBMClassifier	4	3	4	4	5	6	3
AdaBoostClassifier	6	11	6	5	4	10	5
StackingSPINEX	7	6	7	6	13	9	9
SPINEX	9	10	10	7	11	1	4
DecisionTreeClassifier	8	15	8	8	3	5	2
BoostingSPINEX	11	8	12	9	15	14	11
LogisticRegression	15	12	15	11	2	3	1
BaggingSPINEX	10	9	11	12	14	7	8
KNeighborsClassifier	13	13	9	12	1	4	1
SPINEXClassifier(default)	12	14	13	13	12	2	5
SVC	14	7	14	14	7	11	6

672 \*SPINEX = (n\_neighbors=20, distance\_threshold=0.05, distance\_threshold\_decay=0.95, ensemble\_method=None, metric='manhattan')

### 673 3.3 Example with explainability

674 Now, we show one example of the self-interpretability methods included within SPINEX. A sample  
 675 case of synthetic dataset of 500 sample points and 5 features is presented herein. The results of the  
 676 comparison in terms of accuracy and total time/energy used are shown in Fig. 6. As one can see,  
 677 the SPINEX performance is well positioned against the other models.

678 In terms of model explainability, Fig. 7 shows the calculated feature importance as described in  
 679 Sec. 2.0. Comparatively speaking, the calculated trends of feature importance values seem to  
 680 parallel that obtained from other models. The same figure also shows plots of other visualizations  
 681 that can explain model behaviour. Such visualization includes average interaction effects and  
 682 feature combination between features. In addition, the pairwise interactions and feature importance  
 683 at the global level and local level (for a particular instance) are also plotted – please refer to Sec.  
 684 2 for a detailed description of each of these visualizations.

### 685 4.0 Conclusions

686 The SPINEX algorithm offers a novel approach for interpretable regression analysis by integrating  
 687 ensemble learning with feature interaction analysis. It provides accurate predictions while  
 688 unraveling the complex relationships between features and the target variable. The algorithm's  
 689 neighbor-based feature importance and interaction effects offer transparent explanations for  
 690 individual predictions, allowing users to gain insights into the model's decision-making process. It  
 691 is expected that the performance of SPINEX will improve with further development efforts.

692



Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

### 693 **Data Availability**

694 Some or all data, models, or code that support the findings of this study are available from the  
695 corresponding author upon reasonable request.

696 SPINEX and its derivatives can be accessed from [www.mznaser.com](http://www.mznaser.com),  
697 <https://pypi.org/project/SPINEX/> and <https://github.com/mznaser-clemson/SPINEX>.

698 SPINEX can be installed as:

```
699 pip install SPINEX  
700 from SPINEX import SPINEXRegressor  
701 from SPINEX import SPINEXClassifier
```

### 702 **Conflict of Interest**

703 The authors declare no conflict of interest.

### 704 **References**

- 705 [1] J. Too, G. Liang, H. Chen, Memory-based Harris hawk optimization with learning agents:  
706 a feature selection approach, *Eng. Comput.* (2022). <https://doi.org/10.1007/s00366-021-01479-4>.  
707
- 708 [2] I. Naruei, F. Keynia, Wild horse optimizer: a new meta-heuristic algorithm for solving  
709 engineering optimization problems, *Eng. Comput.* (2022). <https://doi.org/10.1007/s00366-021-01438-z>.  
710
- 711 [3] C. Rudin, Stop explaining black box machine learning models for high stakes decisions and  
712 use interpretable models instead, *Nat. Mach. Intell.* (2019). <https://doi.org/10.1038/s42256-019-0048-x>.  
713
- 714 [4] W.J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, B. Yu, Definitions, methods, and  
715 applications in interpretable machine learning, *Proc. Natl. Acad. Sci. U. S. A.* (2019).  
716 <https://doi.org/10.1073/pnas.1900654116>.
- 717 [5] S.N. Van Der Veer, L. Riste, S. Cheraghi-Sohi, D.L. Phipps, M.P. Tully, K. Bozentko, S.  
718 Atwood, A. Hubbard, C. Wiper, M. Oswald, N. Peek, Trading off accuracy and  
719 explainability in AI decision-making: findings from 2 citizens' juries, *J. Am. Med.  
720 Informatics Assoc.* (2021). <https://doi.org/10.1093/jamia/ocab127>.
- 721 [6] H. Ding, I. Takigawa, H. Mamitsuka, S. Zhu, Similarity-based machine learning methods  
722 for predicting drug-target interactions: A brief review, *Brief. Bioinform.* (2013).  
723 <https://doi.org/10.1093/bib/bbt056>.
- 724 [7] G. Dudek, P. Pełka, Pattern similarity-based machine learning methods for mid-term load  
725 forecasting: A comparative study, *Appl. Soft Comput.* (2021).  
726 <https://doi.org/10.1016/j.asoc.2021.107223>.
- 727 [8] T. Hofmann, Learning the similarity of documents: An information-geometric approach to  
728 document retrieval and categorization, in: *Adv. Neural Inf. Process. Syst.*, 2000.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 729 [9] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, Y. Wu, Learning  
730 fine-grained image similarity with deep ranking, in: Proc. IEEE Comput. Soc. Conf.  
731 Comput. Vis. Pattern Recognit., 2014. <https://doi.org/10.1109/CVPR.2014.180>.
- 732 [10] A. Charfi, S. Ammar Bouhamed, E. Bosse, I. Khanfir Kallel, W. Bouchaala, B. Solaiman,  
733 N. Derbel, Possibilistic Similarity Measures for Data Science and Machine Learning  
734 Applications, *IEEE Access*. (2020). <https://doi.org/10.1109/ACCESS.2020.2979553>.
- 735 [11] T. Widiyaningtyas, I. Hidayah, T.B. Adji, User profile correlation-based similarity  
736 (UPCSim) algorithm in movie recommendation system, *J. Big Data*. (2021).  
737 <https://doi.org/10.1186/s40537-021-00425-x>.
- 738 [12] F. Fkih, Similarity measures for Collaborative Filtering-based Recommender Systems:  
739 Review and experimental comparison, *J. King Saud Univ. - Comput. Inf. Sci.* (2021).  
740 <https://doi.org/10.1016/j.jksuci.2021.09.014>.
- 741 [13] P. Domingos, A few useful things to know about machine learning, *Commun. ACM*. (2012).  
742 <https://doi.org/10.1145/2347736.2347755>.
- 743 [14] K. Taunk, S. De, S. Verma, A. Swetapadma, A brief review of nearest neighbor algorithm  
744 for learning and classification, in: 2019 Int. Conf. Intell. Comput. Control Syst. ICCS 2019,  
745 2019. <https://doi.org/10.1109/ICCS45141.2019.9065747>.
- 746 [15] S. Dhanabal, S. Chandramathi, A Review of various k-Nearest Neighbor Query Processing  
747 Techniques, *Int. J. Comput. Appl.* (2011).
- 748 [16] J. Laaksonen, E. Oja, Classification with learning k-nearest neighbors, in: *IEEE Int. Conf.*  
749 *Neural Networks - Conf. Proc.*, 1996. <https://doi.org/10.1109/icnn.1996.549118>.
- 750 [17] D. Bera, R. Pratap, B.D. Verma, Dimensionality Reduction for Categorical Data, *IEEE*  
751 *Trans. Knowl. Data Eng.* (2021). <https://doi.org/10.1109/TKDE.2021.3132373>.
- 752 [18] E.Y. Boateng, J. Otoo, D.A. Abaye, Basic Tenets of Classification Algorithms K-Nearest-  
753 Neighbor, Support Vector Machine, Random Forest and Neural Network: A Review, *J. Data*  
754 *Anal. Inf. Process.* (2020). <https://doi.org/10.4236/jdaip.2020.84020>.
- 755 [19] H.A. Abu Alfeilat, A.B.A. Hassanat, O. Lasassmeh, A.S. Tarawneh, M.B. Alhasanat, H.S.  
756 Eyal Salman, V.B.S. Prasath, Effects of Distance Measure Choice on K-Nearest Neighbor  
757 Classifier Performance: A Review, *Big Data*. (2019).  
758 <https://doi.org/10.1089/big.2018.0175>.
- 759 [20] M.Z. Naser, *Machine Learning for Civil and Environmental Engineers: A Practical*  
760 *Approach to Data-Driven Analysis, Explainability, and Causality*, Wiley, New Jersey, 2023.
- 761 [21] Y. Guo, C. Du, Y. Zhao, T.F. Ting, T.A. Rothfus, Two-level K-nearest neighbors approach  
762 for invasive plants detection and classification, *Appl. Soft Comput.* (2021).  
763 <https://doi.org/10.1016/j.asoc.2021.107523>.
- 764 [22] S.G.K. Patro, B.K. Mishra, S.K. Panda, R. Kumar, H.V. Long, D. Taniar, I. Priyadarshini,

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 765 A Hybrid Action-Related K-Nearest Neighbour (HAR-KNN) Approach for  
766 Recommendation Systems, IEEE Access. (2020).  
767 <https://doi.org/10.1109/ACCESS.2020.2994056>.
- 768 [23] Y.M. Wazery, E. Saber, E.H. Houssein, A.A. Ali, E. Amer, An Efficient Slime Mould  
769 Algorithm Combined with K-Nearest Neighbor for Medical Classification Tasks, IEEE  
770 Access. (2021). <https://doi.org/10.1109/ACCESS.2021.3105485>.
- 771 [24] X. Fan, Z. Chen, L. Zhu, Z. Liao, B. Fu, A Novel Hybrid Similarity Calculation Model, Sci.  
772 Program. (2017). <https://doi.org/10.1155/2017/4379141>.
- 773 [25] J. Gou, L. Du, Y. Zhang, T. Xiong, A new distance-weighted k-nearest neighbor classifier,  
774 J. Inf. Comput. Sci. (2012).
- 775 [26] M. Muja, D.G. Lowe, Scalable nearest neighbor algorithms for high dimensional data, IEEE  
776 Trans. Pattern Anal. Mach. Intell. (2014). <https://doi.org/10.1109/TPAMI.2014.2321376>.
- 777 [27] M. Aumüller, E. Bernhardsson, A. Faithfull, ANN-Benchmarks: A benchmarking tool for  
778 approximate nearest neighbor algorithms, Inf. Syst. (2020).  
779 <https://doi.org/10.1016/j.is.2019.02.006>.
- 780 [28] Y. Dong, X. Ma, T. Fu, Electrical load forecasting: A deep learning approach based on K-  
781 nearest neighbors, Appl. Soft Comput. (2021). <https://doi.org/10.1016/j.asoc.2020.106900>.
- 782 [29] W. Jiang, Time series classification: nearest neighbor versus deep learning models, SN  
783 Appl. Sci. (2020). <https://doi.org/10.1007/s42452-020-2506-9>.
- 784 [30] Z. Dang, C. Deng, X. Yang, K. Wei, H. Huang, Nearest neighbor matching for deep  
785 clustering, in: Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit., 2021.  
786 <https://doi.org/10.1109/CVPR46437.2021.01348>.
- 787 [31] L. Ruff, J.R. Kauffmann, R.A. Vandermeulen, G. Montavon, W. Samek, M. Kloft, T.G.  
788 Dietterich, K.R. Muller, A Unifying Review of Deep and Shallow Anomaly Detection,  
789 Proc. IEEE. (2021). <https://doi.org/10.1109/JPROC.2021.3052449>.
- 790 [32] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in  
791 vector space, in: 1st Int. Conf. Learn. Represent. ICLR 2013 - Work. Track Proc., 2013.
- 792 [33] M. Loog, Nearest neighbor-based importance weighting, in: IEEE Int. Work. Mach. Learn.  
793 Signal Process. MLSP, 2012. <https://doi.org/10.1109/MLSP.2012.6349714>.
- 794 [34] M.Z. Naser, Do We Need Exotic Models? Engineering Metrics to Enable Green Machine  
795 Learning from Tackling Accuracy-Energy Trade-offs, J. Clean. Prod. 382 (2023) 135334.  
796 <https://doi.org/10.1016/J.JCLEPRO.2022.135334>.
- 797 [35] M.Z. Naser, V. Kodur, H.-T. Thai, R. Hawileh, J. Abdalla, V. V. Degtyarev, StructuresNet  
798 and FireNet: Benchmarking databases and machine learning algorithms in structural and  
799 fire engineering domains, J. Build. Eng. (2021) 102977.  
800 <https://doi.org/10.1016/J.JOBE.2021.102977>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 801 [36] M. van Smeden, K.G. Moons, J.A. de Groot, G.S. Collins, D.G. Altman, M.J. Eijkemans,  
802 J.B. Reitsma, Sample size for binary logistic prediction models: Beyond events per variable  
803 criteria:, <https://doi.org/10.1177/0962280218784726>. 28 (2018) 2455–2474.  
804 <https://doi.org/10.1177/0962280218784726>.
- 805 [37] R.D. Riley, K.I.E. Snell, J. Ensor, D.L. Burke, F.E. Harrell, K.G.M. Moons, G.S. Collins,  
806 Minimum sample size for developing a multivariable prediction model: PART II - binary  
807 and time-to-event outcomes, *Stat. Med.* (2019). <https://doi.org/10.1002/sim.7992>.
- 808 [38] I. Frank, R. Todeschini, The data analysis handbook, 1994.  
809 [https://books.google.com/books?hl=en&lr=&id=SXEpb0H6L3YC&oi=fnd&pg=PP1&ots](https://books.google.com/books?hl=en&lr=&id=SXEpb0H6L3YC&oi=fnd&pg=PP1&ots=zfmIRO_XO5&sig=dSX6KJdkuav5zRNxaUdcftGSn2k)  
810 [=zfmIRO\\_XO5&sig=dSX6KJdkuav5zRNxaUdcftGSn2k](https://books.google.com/books?hl=en&lr=&id=SXEpb0H6L3YC&oi=fnd&pg=PP1&ots=zfmIRO_XO5&sig=dSX6KJdkuav5zRNxaUdcftGSn2k) (accessed June 21, 2019).
- 811 [39] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.  
812 Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,  
813 M. Perrot, E. Duchesnay, É. Duchesnay, E. Duchesnay, Scikit-learn: Machine learning in  
814 Python, *J. Mach. Learn. Res.* 12 (2011) 2825–2830.
- 815 [40] R. Kohavi, A study of cross-validation and bootstrap for accuracy estimation and model  
816 selection, *Proc. 14th Int. Jt. Conf. Artif. Intell. - Vol. 2.* (1995).
- 817 [41] T.T. Wong, N.Y. Yang, Dependency Analysis of Accuracy Estimates in k-Fold Cross  
818 Validation, *IEEE Trans. Knowl. Data Eng.* (2017).  
819 <https://doi.org/10.1109/TKDE.2017.2740926>.
- 820 [42] M.Z. Naser, · Amir, H. Alavi, Error Metrics and Performance Fitness Indicators for  
821 Artificial Intelligence and Machine Learning in Engineering and Sciences, *Archit. Struct.*  
822 *Constr.* 2021. 1 (2021) 1–19. <https://doi.org/10.1007/S44150-021-00015-8>.
- 823 [43] A. Khare, Data for Admission in the University, Kaggle. (2022).  
824 [https://www.kaggle.com/datasets/akshaydattatraykhare/data-for-admission-in-the-](https://www.kaggle.com/datasets/akshaydattatraykhare/data-for-admission-in-the-university)  
825 [university](https://www.kaggle.com/datasets/akshaydattatraykhare/data-for-admission-in-the-university).
- 826 [44] M.Z. Naser, V.K. Kodur, Explainable machine learning using real, synthetic and augmented  
827 fire tests to predict fire resistance and spalling of RC columns, *Eng. Struct.* 253 (2022)  
828 113824. <https://doi.org/10.1016/j.engstruct.2021.113824>.
- 829 [45] I.-C.C. Yeh, Modeling of strength of high-performance concrete using artificial neural  
830 networks, *Cem. Concr. Res.* 28 (1998) 1797–1808. [https://doi.org/10.1016/S0008-](https://doi.org/10.1016/S0008-8846(98)00165-3)  
831 [8846\(98\)00165-3](https://doi.org/10.1016/S0008-8846(98)00165-3).
- 832 [46] M.Z. Naser, AI-based cognitive framework for evaluating response of concrete structures  
833 in extreme conditions, *Eng. Appl. Artif. Intell.* 81 (2019) 437–449.  
834 <https://www.sciencedirect.com/science/article/pii/S0952197619300466> (accessed April 1,  
835 2019).
- 836 [47] S. Thai, H.-T. Thai, B. Uy, T. Ngo, M.Z. Naser, Test database on concrete-filled steel  
837 tubular columns, *Mendeley*, 2020. <https://doi.org/10.17632/3XKNB3SDB5.5>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

- 838 [48] U. Chowdhury, Energy Efficiency Data Set, Kaggle. (2022).  
839 <https://www.kaggle.com/datasets/ujjwalchowdhury/energy-efficiency-data-set>.
- 840 [49] Fedesoriano, Body Fat Prediction Dataset, Kaggle2. (2021).  
841 <https://www.kaggle.com/datasets/fedesoriano/body-fat-prediction-dataset>.
- 842 [50] P. Cortez, A. Morais, Forest Fires Data Set Portugal | Kaggle, (2007).  
843 <https://www.kaggle.com/datasets/ishandutta/forest-fires-data-set-portugal> (accessed July  
844 11, 2022).
- 845 [51] Devphaib, Estimating the age of abalone at a seafood farm, Kaggle. (2022).  
846 [https://www.kaggle.com/datasets/devzohaib/estimating-the-age-of-abalone-at-a-seafood-](https://www.kaggle.com/datasets/devzohaib/estimating-the-age-of-abalone-at-a-seafood-farm)  
847 [farm](https://www.kaggle.com/datasets/devzohaib/estimating-the-age-of-abalone-at-a-seafood-farm).
- 848 [52] Fedesoriano, Synchronous Machine Dataset, Kaggle. (2022).  
849 <https://www.kaggle.com/datasets/fedesoriano/synchronous-machine-dataset>.
- 850 [53] R. Patel, RETAIL ANALYSIS WITH WALMART SALES DATA, Kaggle. (2021).  
851 <https://www.kaggle.com/datasets/rutuspatel/retail-analysis-with-walmart-sales-data>.
- 852 [54] M.K. al-Bashiti, M.Z. Naser, Verifying domain knowledge and theories on Fire-induced  
853 spalling of concrete through eXplainable artificial intelligence, *Constr. Build. Mater.* 348  
854 (2022) 128648. <https://doi.org/10.1016/J.CONBUILDMAT.2022.128648>.
- 855 [55] Pima Indians Diabetes Database, Kaggle. (2016).  
856 <https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>.
- 857 [56] M. Abedi, M.Z. Naser, RAI: Rapid, Autonomous and Intelligent machine learning approach  
858 to identify fire-vulnerable bridges, *Appl. Soft Comput.* (2021).  
859 <https://doi.org/10.1016/j.asoc.2021.107896>.
- 860 [57] B.A. Young, A. Hall, L. Pilon, P. Gupta, G. Sant, Can the compressive strength of concrete  
861 be estimated from knowledge of the mixture proportions?: New insights from statistical  
862 analysis and machine learning methods, *Cem. Concr. Res.* 115 (2019) 379–388.  
863 <https://doi.org/10.1016/j.cemconres.2018.09.006>.
- 864 [58] W.H. Wolberg, O.L. Mangasarian, Multisurface method of pattern separation for medical  
865 diagnosis applied to breast cytology, *Proc. Natl. Acad. Sci. U. S. A.* (1990).  
866 <https://doi.org/10.1073/pnas.87.23.9193>.
- 867 [59] W. Wolberg, Breast Cancer Wisconsin (Original) Data Set, UCI Mach. Learn. Repos. (n.d.).  
868 <https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Original%29>.
- 869 [60] M. Koklu, Rice Dataset Commeo and Osmancik, Kaggle. (2022).  
870 <https://www.kaggle.com/datasets/muratkokludataset/rice-dataset-commeo-and-osmancik>.
- 871 [61] R. Saluja, Bank Note Authentication UCI data, Kaggle. (2018).  
872 <https://www.kaggle.com/datasets/ritesaluja/bank-note-authentication-uci-data>.
- 873 [62] A. Kadiwal, Water Quality, Kaggle. (2021).

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.asoc.2024.111518>.

Please cite this paper as:

Naser M.Z., Al-Bashiti M.K., Naser A.Z. (2024). SPINEX: Similarity-based Predictions with Explainable Neighbors Exploration for Regression and Classification. *Applied Soft Computing*. <https://doi.org/10.1016/j.asoc.2024.111518>.

874 <https://www.kaggle.com/datasets/adityakadiwal/water-potability>.

875 [63] S. Bansal, Machine Predictive Maintenance, Kaggle. (2021).

876 <https://www.kaggle.com/datasets/shivamb/machine-predictive-maintenance-classification>.

877 [64] D. Babativa, Depression dataset, Kaggle. (2023).

878 <https://www.kaggle.com/datasets/diegobabativa/depression>.

879 [65] G. Santello, Cars - Purchase Decision Dataset, Kaggle. (2022).

880 <https://www.kaggle.com/datasets/gabrielsantello/cars-purchase-decision-dataset>.

881