

Explainable Artificial Intelligence and Machine Learning in Health Care

Xi Chen

Department of Advanced Computing Sciences

Faculty of Science and Engineering

Maastricht University

Maastricht, The Netherlands

Abstract—Artificial intelligence (AI) holds promise for improving healthcare outcomes, but its adoption in clinical practice is hindered by the lack of transparency in current models. Explainable AI (XAI) addresses this issue by making AI systems more interpretable. This paper explores explainable AI in the healthcare domain, focusing on illness detection. The study evaluates and compares various explainability methods using the Pima Indian Diabetes and Parkinson’s Disease datasets. Both model-agnostic and model-specific methods are explored, including Partial Dependence Plot, Accumulated Local Effects, permutation Feature Importance, SHAP, and model-specific methods for tree structures.

The results reveal that the chosen methods have distinctive strengths and weaknesses, making them more or less suitable for specific applications and datasets. Due to the high expressive power of decision trees, they can be useful in practical procedures. As a result, a human-understandable explanation in the form of a decision tree is generated and optimized based on the study results for both datasets.

I. INTRODUCTION

The field of artificial intelligence (AI) and machine learning have witnessed remarkable advancements in recent years, leading to diverse applications across sectors such as finance, engineering, criminal justice, and healthcare (Carvalho, Pereira, & Cardoso, 2019) [1]. However, a pervasive challenge encountered with these technologies lies in their lack of transparency, leaving end-users without a clear understanding of the decision-making rationale. Instead of providing explanations, AI models typically present prediction results, raising concerns regarding ethical considerations and fairness.

Within the medical and healthcare domain, explainability assumes a critical role, particularly in the context of illness detection. Clinicians should exercise caution when dealing with the outcomes of machine learning techniques, as these outcomes significantly impact consequential decision-making processes where erroneous choices can have severe consequences, including health risks. Establishing trust in healthcare practices is of utmost importance, as it empowers clinicians with a deeper understanding of how the model works, enabling them to effectively identify and rectify errors. This, in turn, fosters patients’ trust in the healthcare system and ensures they receive optimal care, leading to improved health outcomes. Moreover, explainability fosters the discovery of previously unidentified patterns and insights, while simultaneously enhancing the development and utilization of AI in healthcare.

By providing transparency and interpretability, explainable AI methods help overcome regulatory challenges and ethical considerations, promoting the responsible and widespread adoption of AI technologies in the healthcare industry. [2]

As the lack of transparency poses a significant impediment to the effective deployment of AI within the healthcare domain, the emergence of Explainable AI presents promising avenues for addressing this challenge. The seminal work by Markus, Kors, and Rijnbeek (2021) [3] has introduced a comprehensive framework that facilitates the discernment of suitable explainability methods in healthcare, encompassing diverse approaches such as model-based, attribution-based, and example-based explanations, along with global and local perspectives. This framework plays a pivotal role in guiding the design of explainable AI systems within the healthcare domain, formalizing a procedural framework that engenders clarity and comprehensibility. Furthermore, scholarly investigations from various disciplines, including psychology and social perspectives, have delved into the conceptual underpinnings of explainability, unveiling essential properties that inform the evaluation and assessment of explainability methods and explanations (Islam et al., 2019) [4]. These identified properties serve as invaluable guidelines for discerning the utility and efficacy of explainability in the healthcare context.

However, building a trustworthy AI system is a challenging task that requires extensive examination and validation in various real-world cases, especially in the field of illness detection. It is crucial to not only assess its effectiveness but also continuously improve its performance when applied to practical healthcare scenarios. Furthermore, there is a pressing need to prioritize the development of clear and human-understandable explanations that are practical and applicable in healthcare settings. It is important to recognize that explanations for illness diagnosis may differ significantly from explanations in other fields, thus necessitating specific improvements and adjustments tailored to healthcare requirements.

Additionally, it is worth noting that the majority of existing explainability methods primarily focus on generating explanations for supervised learning algorithms, leaving a gap in addressing the needs of unsupervised learning in healthcare. However, unsupervised learning can offer distinct advantages such as improved performance in scenarios where obtaining a complete set of labels may be challenging or unfeasible, as

well as the ability to discover underlying patterns. To overcome this limitation, more efforts should be directed toward developing and adapting unsupervised learning algorithms specifically for illness detection to classification tasks and integrating them with explainability techniques. This holistic approach, encompassing both supervised and unsupervised learning paradigms, would ensure a more comprehensive and accurate understanding of AI models' decision-making processes in healthcare.

A. Problem Statement

This paper aims to review and examine existing explainability methods for certain illness detection by evaluating their performance on two classic datasets: the Parkinson Dataset and the Pima Indian Diabetes Dataset. Additionally, the paper seeks to define evaluation metrics for assessing the quality of explanations generated by these methods. Furthermore, the study aims to adapt unsupervised learning algorithms specifically for illness detection and integrate them with explainability techniques. The ultimate goal is to provide clear and human-understandable explanations that can be practically utilized during the illness-detection decision-making process. The main research questions guiding this study are:

- 1) How to measure the utility of the explanation generated by different methods?
- 2) How do the explainability methods perform given different datasets and what could be the reason if the performance of the methods differs?
- 3) Whether model-specific methods achieve better results than model-agnostic methods?
- 4) Whether model-agnostic methods could be used on the results of unsupervised learning algorithms?
- 5) How to improve the performance of the methods?

Answering these research questions can enhance our understanding of explainability methods in healthcare, specifically in illness detection. It enables us to identify effective approaches, pinpoint areas for improvement, and guide future research.

B. Outline

The structure of the paper will be as follows: Section II to VI are methodology sections, which will clearly define the approach taken: Section II define the general notation for all mathematical formula. Section III presents the datasets used in the study, providing a detailed description of their characteristics. In Section IV, the machine learning algorithms employed in the analysis are discussed. Section V delves into the explainability methods utilized in the research. Then, Section VI outlines the chosen evaluation criteria for assessing the performance of the explainability methods. Section VII will show the implementation of the aforementioned methods and the libraries used. Section VIII will present the experimental results by applying the different methods to the datasets, highlighting and comparing their performance. Section IX will rigorously evaluate the performance of the methods and discuss the findings. Finally, in sections X conclusions will

be drawn, and potential avenues for further research will be suggested.

II. NOTATION

This section introduces the general notation used throughout the thesis to represent various formulas.

An ML model, denoted as \hat{f} , is a prediction function that maps a feature vector from input space \mathcal{X} to a prediction in output space \mathcal{Y} , represented as $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$. The dataset $S = \{z_i\}_{i=1}^n = \{(x_i, y_i)\}_{i=1}^n$ is sampled from a distribution D over a domain $Z = \mathcal{X} \times \mathcal{Y}$. Here, \mathcal{X} represents the instance domain as a set, \mathcal{Y} represents the label domain as a set, and $Z = \mathcal{X} \times \mathcal{Y}$ represents the example domain as a set. n represents the number of data instances.

III. DATASETS

This section provides an overview of the datasets used.

1) *Parkinson Dataset*: The Parkinson's Dataset used in this study is based on the work of Max A. Little et al. (2008) [5]. It includes voice recordings from 195 instances, both with and without Parkinson's disease. The dataset contains voice-related features and diagnostic labels indicating Parkinson's disease presence. Limitations of the dataset include limited variability due to sampling from a small size of 31 individuals (23 with PD) and the size of the dataset is relatively small, which may impact model generalization and pose challenges in training and evaluation.

2) *Diabetes Dataset*: The Pima Indian Diabetes Dataset, introduced by Smith et al. (1988) [6], comprises data from 768 Pima Indian women in Phoenix, Arizona, USA. It includes a diverse set of human body features and diagnostic labels indicating the presence or absence of diabetes. The dataset exhibits appropriate variance given its size. However, the dataset presents challenges due to its strong interactions and non-linear relationships among the features.

IV. MACHINE LEARNING ALGORITHMS

This section provides an overview of the machine learning algorithms utilized in the study for the classification task of illness detection.

A. Supervised Algorithms

The first algorithm used is the **Decision Tree**, which constructs a tree-like model to classify instances based on if-then conditions. Decision Trees are highly interpretable and offer insights into the decision-making process. They serve as a reference for Proxy Model Comparison in assessing fidelity properties.

A **Random Forest** ensemble is employed, combining multiple decision trees for predictions. While individual trees may be interpretable, the overall model's interpretability is reduced.

KNN (k-nearest neighbors) is a non-parametric method that assigns class labels based on similarity to neighbors. Interpreting KNN can be intricate, but the concept of using neighbors for classification is comprehensible.

Multilayer Perceptron (MLP) is also included in this study. It is an artificial neural network with interconnected

layers. MLPs generally have lower interpretability due to complex weight and bias relationships.

Support Vector Machines (SVMs) are used, finding optimal hyperplanes for class separation using the kernel trick. Non-linear SVMs, especially with high-dimensional spaces, pose interpretation challenges. While the hyperplane may be interpretable, understanding individual feature relationships becomes complex. [7]

B. Unsupervised Algorithms

Unsupervised learning algorithms are not directly used for classification tasks, but they can be transformed into semi-supervised learning [8]. In this context, the goal is to assign observations to predefined classes using both labeled and unlabeled data. By incorporating class information, these semi-supervised clustering methods bridge the gap between unsupervised and supervised learning, enabling the use of unsupervised techniques for classification.

To illustrate this approach, one popular clustering algorithm **K-means** is selected. It effectively partitions data into distinct groups based on proximity to cluster centroids. The algorithm iteratively assigns data points to the nearest centroid and updates centroids by calculating the mean of assigned points until convergence. To implement this approach, the K-means algorithm is applied to the dataset, associating each cluster with the majority label of its instances. Instances are then assigned the label of the closest centroid. However, in this report, only labeled data was used and tested.

Although not strictly unsupervised, this approach offers advantages similar to unsupervised methods. It handles missing data effectively, which is crucial for real-time illness diagnosis scenarios. The method demonstrates robustness to missing labels by relying on the majority within each cluster. It is also computationally efficient, enabling the classification of large datasets. Furthermore, it can uncover underlying patterns and structures in the data, enhancing its usefulness in illness detection tasks. [9]

V. EXPLAINABILITY METHODS

The following explainability methods are employed on the classifiers discussed in Section IV to analyze their behavior. The purpose is to evaluate the efficacy of these explainability methods and their corresponding explanations. Furthermore, their application is explored to enhance the utilization of AI models in healthcare, specifically in illness detection. [11] [10] [12]

A. Partial Dependence Plots (PDP)

Partial Dependence Plots (PDPs) visualize how the predicted outcome changes with variations in selected features while keeping other features constant. They illustrate the impact of features on the probability of the positive class in binary classification, providing attribution-based explanations.

The partial dependence function in PDP is estimated by manipulating the selected feature while holding other features constant and averaging the predictions across all data

instances. Equation 1 represents the basic formula for calculating the marginalized classifier $\hat{f}_{x_S}(x_S)$, which considers the selected features x_S and their complement $x_C^{(i)}$, with $\mathcal{X} = x_S \times x_C^{(i)}$.

$$\hat{f}_{x_S}(x_S) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x_S, x_C^{(i)}) \quad (1)$$

Despite being useful, PDPs have limitations. They may not accurately represent correlated features and struggle to capture nonlinear relationships, potentially misrepresenting feature relationships and interactions. To overcome these limitations, alternative methods like ALE plots are recommended. It's crucial to interpret PDPs cautiously and avoid oversimplified conclusions about feature importance and model behavior.

B. Accumulated Local Effects (ALE)

The Accumulated Local Effects (ALE) plots capture the accumulated average marginal effects of a specific feature on the model's predictions. They show how the average prediction changes as the feature's value varies across the dataset. In binary classification, ALE reveals the average effect of a feature on the predicted probability of the positive class, considering other feature effects. ALE generates attribution-based explanations, as it shows the average effect of input features on the output. The basic ALE formula can be illustrated by :

$$\hat{f}_{S,ALE}(x_S) = \int_{z_{0,S}}^{x_S} E_{X_C|X_S=x_S} [\hat{f}^S(X_s, X_c)|X_S = z_S] dz_S - c \quad (2)$$

The equation 2 represents the definition of ALE for the selected feature x_S . In this equation, $\hat{f}^S(x_S, x_C)$ represents the local effects of the feature x_S on the function $\hat{f}(\cdot)$ at the specific point (x_1, x_S) . The integral from $z_{0,S}$ to x_S captures the accumulated average of these local effects, considering different values of x_S . The conditional expectation $E_{X_C|X_S=x_S} [\hat{f}^S(X_s, X_c)|X_S = z_S]$ represents the expected value of the local effects $\hat{f}^S(X_s, X_c)$ given $X_S = z_S$, where X_C follows a conditional density $P(X_C|X_S)$. This accounts for the varying influence of x_S at different values of x_S . To establish a reference point, a constant c is subtracted from the integral, which provides a baseline for comparison.

ALE plots excel in analyzing correlated features, accurately depicting their relationship with the predicted outcome. Unlike PDPs, ALE plots can capture complex interactions and dependencies among correlated features. However, they have limitations in representing nonlinear relationships, which may restrict their ability to capture intricate feature interactions.

C. Permutation Feature Importance

Permutation Feature Importance assesses the significance of individual features by shuffling their values and measuring the resulting change in model performance. Thus, it is an attribution-based explanation generator. The basic Permutation Feature Importance formula can be illustrated as follows:

$$FI^{perm} = \frac{1}{N} \sum_{i=1}^N (e_{perm} - e_{orig}) \quad (3)$$

$$e_{orig} = L(y, \hat{f}(X)), e_{perm} = L(y, \hat{f}(X_{perm})) \quad (4)$$

The importance FI of a feature x_S is determined by comparing the model's performance before permutation e_{orig} with the performance after permutation e_{perm} , where e denotes the estimated loss. The magnitude of the decrease in model performance provides insights into the feature's contribution, highlighting whether the model heavily relies on the information provided by that feature for accurate predictions.

In the formula, N represents the number of samples used for the calculation. By averaging the differences between e_{perm} and e_{orig} overall N samples, the Permutation Feature Importance is computed. This explainability method is effective for identifying influential features, regardless of their linear or nonlinear relationships with the predicted outcome. It handles correlated features well and captures complex interactions among them. However, it has limitations because it assumes feature independence and the ability to arbitrarily shuffle their values, which may not always be the case.

D. SHapley Additive exPlanations (SHAP)

SHapley Additive exPlanations (SHAP) is a versatile framework that explains predictive models by assigning values to each feature based on their contribution to the predictions. It offers both global and individual explanations and is based on Shapley values from cooperative game theory.

The basic SHAP methodology provides explanations using a structured format, where the explanatory model g is applied to a coalition vector $z' \in \{0, 1\}^M$ constructed from the sample data instances z . Here, M represents the total number of features used in the data instances. The Shapley value for feature j is denoted as ϕ_j . For explaining a specific instance x , the coalition vector x' consists of all ones, resulting in a simplified explanation model as demonstrated by:

$$g(z') = \phi_0 + \sum_{j=1}^M \phi_j z'_j \quad (5)$$

$$g(x') = \phi_0 + \sum_{j=1}^M \phi_j \quad (6)$$

Kernel SHAP and Tree SHAP are utilized to estimate Shapley values in this research, explaining feature contributions to model predictions. Kernel SHAP combines predictions through weighted averaging, while Tree SHAP uses a tree-based model to partition the feature space and assign prediction differences.

SHAP provides both global (Attribution-based) and individual (Example-based explanations), enabling a comprehensive understanding of feature importance and its influence on predictions. However, it should be noted that SHAP does not establish causal relationships.

E. Decision Trees as Model-Based Explanations

Decision trees are extensively utilized in machine learning to facilitate explainability and fall under the category of model-based explanations. Within the realm of healthcare, decision

trees serve a dual purpose as both explainable modeling techniques and post-hoc explanations.

Explainable modeling entails the development of a task model that possesses inherent interpretability, catering specifically to user groups such as clinicians. Decision trees are highly regarded for their interpretability, primarily owing to their transparent structure comprising a sequence of binary decisions based on input features. Each decision node represents a distinct feature, and the paths traversed from the root to the leaf nodes dictate the predicted outcome. Consequently, clinicians can comprehensively comprehend and track the decision-making process of the model. In the context of post-hoc explanations, decision trees can be constructed as supplementary models to shed light on the underlying rationale of the task model. This approach offers a more interpretable alternative to complex models, enabling clinicians to gain valuable insights into the influential factors influencing the predictions. However, it is crucial to acknowledge that post-hoc explanations may not fully capture the intricacies of the original task model, necessitating careful consideration during practical sessions. [3]

F. Tree Based Explainability Method

Gini Importance measures the importance of features in tree-based models, such as random forests and decision trees. It quantifies the contribution of each feature by evaluating the reduction in impurity achieved through its inclusion in the model. Higher impurity reduction indicates greater feature importance.

Tree visualization and decision rules are model-specific methods used to interpret decision trees. Tree visualization provides a graphical representation of the tree structure, aiding in understanding the hierarchy of decisions. Decision rules extract concise and logical statements that describe the conditions leading to specific predictions.

VI. EVALUATION CRITERIA

The application of explainability methods in healthcare, specifically in illness detection, is an important and evolving research area. However, the effectiveness of these methods can vary depending on the dataset and model used. To ensure valuable explanations, it is necessary to evaluate their performance using defined criteria. Establishing performance criteria allows for structured assessment and comparison of different techniques. By applying these criteria to datasets with specific characteristics, researchers can examine the impact of challenges on explainability methods. While domain specificity limits rating assignment in this research, using performance criteria enables objective evaluation of explainability methods in illness detection. Evaluation of an explainability method depends on both the method itself and the quality of generated explanations. Effectiveness requires assessing the method and its explanations in facilitating human understanding of the machine learning model. While most evaluation criteria are from reference [4] [12], modifications were made, including the removal of some criteria and the addition of new ones

by myself like generalizability and scope. Comprehensive practical evaluation methods have been developed to assess these properties along with critical defined properties.

A. Evaluation Criteria for explainability methods

- 1) Expressive Power: Focuses on the method's ability to effectively capture and convey the underlying reasoning of a machine learning model. It includes various representations like decision trees, IF-THEN rules, weighted sums, decision graphs, and natural language.
- 2) Translucency: Refers to how much the explanation method relies on accessing the internal mechanisms of the model. Transparent models, like decision trees, exhibit high translucency, while model-agnostic methods possess zero translucency. Higher translucency involves utilizing model-specific information for more specific explanations.
- 3) Portability: Indicates the range of machine learning models that the explanation method can cover effectively. Methods with lower translucency, treating the model as a black box, tend to have greater portability as they don't rely on model-specific information. However, there's a trade-off between translucency and portability.
- 4) Algorithmic Complexity: Measures the computational time required to generate explanations. Some methods may be more computationally demanding, impacting their usability, especially with large datasets or real-time scenarios.

B. Evaluation Criteria for explanations

- 1) Fidelity: Refers to the extent to which an explanation approximates the predictions of the black-box model. It will be evaluated based on Proxy Model Comparison: Training an interpretable model, like a decision tree, to mimic the black-box model's behavior. Comparing the explanations generated by the surrogate model with the black-box model's explanations. Similarity indicates higher fidelity.
- 2) Consistency: Measures the similarity among explanations generated for different machine learning models for the same task or dataset.
- 3) Stability: Examines the similarity of explanations for similar instances within a specific machine learning model, considering the consistency of explanations with slight variations in input data.
- 4) Generalizability: Evaluates the performance of explanation methods across different datasets to assess their effectiveness in various contexts.
- 5) Comprehensibility: Reflects the extent to which the recipient of the explanation can understand it. Measures may include the number of features with non-zero weights or rules in a decision tree. A quality analysis based on personal experience is conducted in the context of this report.
- 6) Degree of Importance: Evaluates how well an explanation covers essential features or reflects key aspects,

considering the placement of rules or features within the explanation.

- 7) Scope: Assesses whether the explainability method provides global or local explanations, capturing the overall behavior of the model or focusing on specific instances or subsets of data.

VII. IMPLEMENTATION

This sections show the implementation of the aforementioned methods and the libraries used.

A. Machine Learning Algorithm

The supervised machine learning algorithms mentioned in Section 3, including decision tree, random forest, K-nearest neighbors (KNN), multi-layer perceptron (MLP), and support vector machine (SVM), are implemented using the scikit-learn library [13]. Careful parameter tuning is performed for each classifier to optimize their performance based on the specific characteristics and complexities of the given datasets. Parameters such as the number of iterations or tree depth are adjusted to achieve the best results.

Furthermore, an additional classification method based on an unsupervised algorithm is incorporated. This method employs the K-means clustering algorithm as its foundation. Each cluster formed by K-means is then associated with the majority label of the instances it contains. To classify the test dataset, the closest centroid is determined for each test instance, and the corresponding label is assigned accordingly. This unsupervised approach leverages the clustering structure to assign labels based on the majority within each cluster.

B. Explainability Method

The explainability methods utilized in this study make use of various libraries and code implementations for their application. Specifically:

The scikit-learn library [13] is suitable for Gini Importance and decision tree analysis. The SHAP library [15] enables SHapley Additive exPlanations for global and individual insights. For Permutation Feature Importance, the 'permutation_importance' function in sklearn [13] is useful. The 'PDPbox' library [16] facilitates Partial Dependence Plots (PDPs), while the 'PyAle' package [14] is suitable for ALE plots.

It should be noted that the mathematical formula provided in the Explainability Method Section is slightly different in format from the actual implementation in the library due to practical considerations. However, it is important to emphasize that they are fundamentally the same in theory and serve the same purpose of capturing and interpreting the feature importance and effects.

VIII. EXPERIMENTS AND RESULTS

This section presents the experimental setup and the corresponding results obtained from applying various explainability Methods to the machine learning algorithms.

A. Experimental Setup

The experiment employed an 80% training and 20% test data split to train machine learning models on two datasets, with a fixed random seed of 42. The k-means-based classifier was configured with 7 clusters for the Parkinson's dataset and 9 clusters for the Diabetes dataset. A subset of 50 data points was selected for the ALE and PDP analyses. The accuracies of the models are presented in Appendix 10, and Appendix 7 provides an overview of the explainability methods used for each model. Additional details of the experimental settings can be found in Appendix.

B. Parkinson Dataset

1) *PDP & ALE:* This section demonstrates the explanatory analyses conducted on the Parkinson dataset using PDP and ALE. These methods both offer insights into the relationship between features and positive class prediction. Each figure represents a distinct feature.

Out of the 132 groups of ALE and PDP results, only 38 groups show consistent trends. Among these groups, important features demonstrate variations in their impact as their values change (6 groups). This suggests that these features play a significant role in the prediction. For example, Figure 1 displays the PDP and ALE explanations for the "MDVP:Fo(Hz)" feature. The plots indicate a substantial decrease around 200, followed by a stable trend, indicating that values exceeding 200 are associated with a lower likelihood of Parkinson's disease diagnosis. In contrast, unimportant features consistently exhibit ALE and PDP values of 0 (32 groups).

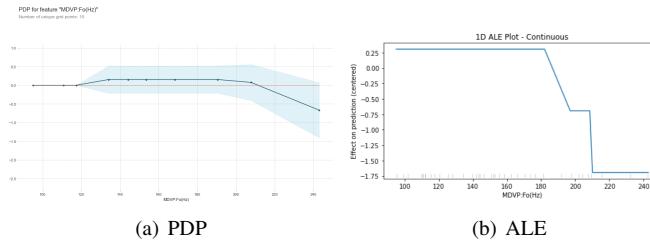


Fig. 1. PDP and ALE Explanations for the "MDVP:Fo(Hz)" Feature in the Decision Tree

When comparing the ALE and PDP results, a common observation (85 out of 132 cases) is that PDP tends to generate a curve that closely follows the zero line, while ALE produces a straight line at ale=0. An example illustrating this difference can be seen in Figure 2, which demonstrates the impact of the "MDVP:Fo(Hz)" feature on the random forest model. This phenomenon can be attributed to the characteristics of the dataset, which exhibits low variability, as well as the inherent characteristics of the machine learning algorithms used. Further discussion on these aspects will be provided later. Furthermore, a subset of features (9 groups out of 132) demonstrates distinct patterns, suggesting the presence of non-linear or interactive effects. This observation is exemplified in Figure 3, which showcases the effect of the "MDVP:Fo(Hz)" feature in the MLP model.

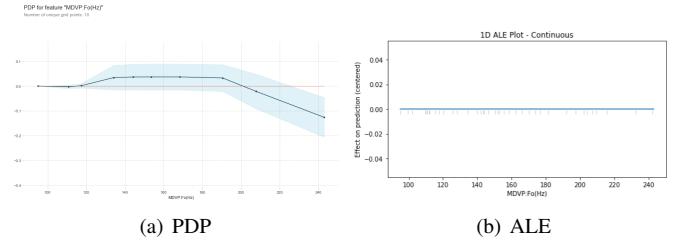


Fig. 2. PDP and ALE Explanations for the "MDVP:Fo(Hz)" Feature in the random forest

For a more comprehensive analysis and statistical comparison between PDP and ALE, please refer to Appendix. Additionally, all the PDP and ALE figures can also be found in Appendix.

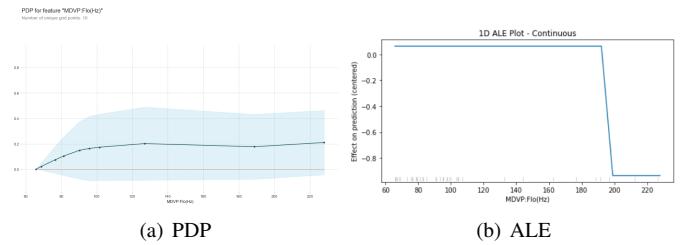


Fig. 3. PDP and ALE Explanations for the "MDVP:Fo(Hz)" Feature in the MLP

Additionally, it is noteworthy that the majority of PDP and ALE plots consistently exhibit values predominantly above the horizontal line at 0. Further analysis and discussion on this observation will be provided in subsequent sections.

2) *SHAP(global), Permutation Importance, & Gini Importance:* This section focuses on the explanations provided by the SHAP, Permutation Importance, and Gini Importance methods, which offer feature importance rankings and numerical weights.

As exemplified in Figure 4, the explanation from a decision tree model is presented in a visually intuitive manner, with bar lengths representing the importance of each feature. The influential features, such as "PPE" and "MDVP:FO(Hz)," receive the highest importance rankings. They are closely followed by "NHR," "MDVP:RAP," "spread1," "D2," and "RPDE."

In this context, influential features are defined as those with importance scores greater than the non-zero median feature importance score. Within the same ML model, it is observed that the ranking of important features remains highly similar, especially for the top few features, indicating stability and consistency. These methods tend to align with each other to some extent, although there may be slight differences in the rankings. Regarding the similarity among explanations from different models, although the rankings may vary, there are certain important features that are consistently considered significant. For example, the feature "MDVP:FO(Hz)" is deemed important by 100% of the ML models, while "spread1" and "PPE" are regarded as important by 4 out of 6 models.

For a more detailed overview of the influential features tables, including a comparison of important features for different combinations of ML models and explainability methods, please refer to Appendix. Further discussion of these findings will be provided later.

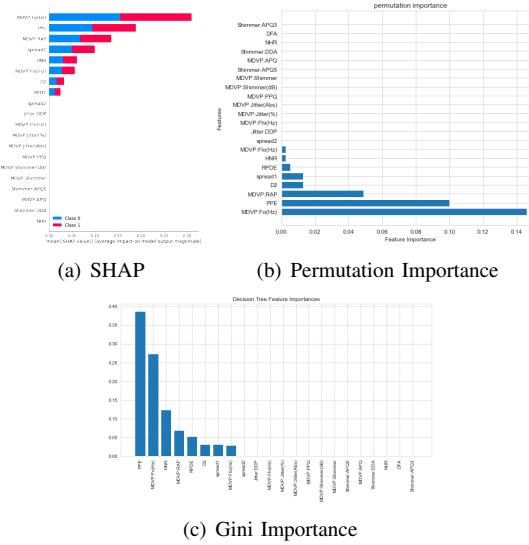


Fig. 4. Various Feature Importance Methods in the Decision Tree

Additionally, Longer bars in SHAP plots indicate higher impact, whether positive or negative, on the model's prediction. Red bars represent a positive impact, while blue bars represent the opposite. Negative permutation importances can occur when shuffled or noisy data produces better predictions than the original data, indicating minimal feature significance ideally close to zero. However, random chance may lead to more accurate predictions on shuffled data, particularly in smaller datasets where chance plays a larger role due to limited data.

3) Example-based Explanation: SHAP (Local) & Decision Tree: SHAP and Decision Tree are the only two explainability methods that provide explanations at an individual level. Decision graphs serve as an effective visualization technique for SHAP individual explanations, showcasing the hierarchical influence of specific features on predictions. On the other hand, decision trees generate explanations based on the decision path, resulting in potentially different explanations for different instances.

Figure 5 presents the SHAP plot for two specific instances, namely no. 113 and 144 (details provided in Appendix), utilizing a decision tree model. These instances are highly similar, with minimal differences of less than 5% for all influential features. However, their predictions are opposite, as they lie at the boundary of the decision tree's one-end branch. It can be observed that both instances have the same feature importance ranking and weights. However, the instance with the false prediction follows a left path in the decision graph, while the positive prediction to the right. This difference is due to the final decision made by the decision tree, leading to

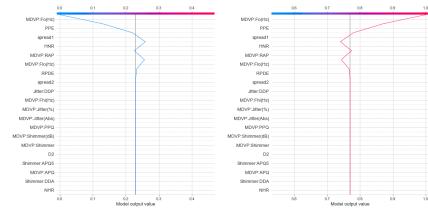


Fig. 5. Two Similar Sample's Decision Graphs in the Decision Tree

different outcomes. This variation also affects the calculation of SHAP values, where features with positive SHAP values have influenced the model to make a positive prediction. The same principle applies to the decision tree's last step in the decision path.

However, when using Decision Graphs with other more robust ML models (e.g., random forest, SVM) or instances that are not close to the boundary of the decision tree, the outcomes are more consistent. For further details, please refer to the appendices.

Additionally, the feature weights in individual explanations mostly correspond to the feature importance derived from the global SHAP analysis. This consistency holds true for correctly predicted instances as well as falsely predicted ones. It is observed that individual explanations from SHAP align with global explanations.

4) Decision Tree, Decision Rules, and Tree Visualization:

As discussed, The decision tree serves as a model-based explanation due to its white-box nature, combining decision rules, tree visualization, and Gini importance. Given the low complexity of the dataset, the decision tree techniques result in manageable sizes and easy-to-read representations. Detailed visualizations of the decision tree techniques can be found in Appendix.

C. Diabetes Dataset

1) PDP & ALE: As the Diabetes Dataset showcases proper variability, the ALE plots do not consistently display a straight line at 0, unlike in the Parkinson dataset. However, in the case of the diabetes dataset, 27 out of the 48 groups of PDP and ALE results exhibit distinct patterns. Figure 6 provides an example of such a distinct pattern utilizing MLP. One possible explanation for this is the presence of strong interactive effects within the dataset, such as the relationship between the number of pregnancies and age. However, this is not the sole reason, as ALE is also effective at capturing relationships when they exist. However, the variations in ALE results across different ML models for the same feature suggest the bad performance for ALE and the possible presence of nonlinear relationships.

2) SHAP(global), Permutation Importance, & Gini Importance: Similar to the findings in the Parkinson dataset, the feature importance generated for the Diabetes dataset demonstrates a remarkable similarity within the same ML model, and even across all ML models. Among all the models, the

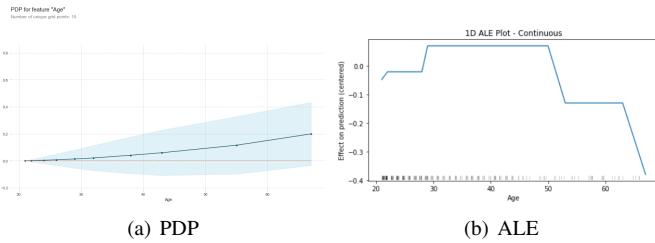


Fig. 6. PDP and ALE Explanations for the "Age" Feature in the MLP

feature "Glucose" consistently emerges as the most important, with 5 out of 6 models considering it the top influential feature. Additionally, "BMI" and "Age" are also identified as influential features in 4 out of 6 models.

3) *Example-based Explanation: SHAP (Local) & Decision Tree:* As observed in the Parkinson's dataset, similar findings apply in this context as well. Decision trees are sensitive when instances are near the decision boundary, and even slight changes can lead to different explanations and final results. This highlights the importance of addressing this issue in optimization when applying decision trees in practical scenarios.

Regarding SHAP, additional samples can be found in Appendix. Similar to the Parkinson's dataset, instances that are similar to each other tend to generate more consistent explanations in SVM and random forest models.

4) *Decision Tree, Decision Rules, and Tree Visualization:* Due to the increased complexity of the diabetes dataset, the size of the decision tree is larger, which somewhat diminishes its expressive power. Detailed visualizations can be found in Appendix.

D. Model Optimization in the Context of Illness Detection

In illness detection, accuracy plays a critical role in health-care outcomes. Decision tree models are preferred for their interpretability and practicality in clinical procedures. However, to enhance performance and address uncertainty near decision boundaries, an optimization approach is proposed. It involves minimizing false predictions by assigning the "unknown" class to cases in proximity to the decision boundary. This reduces the risk of misclassification and mitigates false predictions. Individuals labeled as "unknown" are advised to seek professional help. The approach requires setting a threshold to define the range of feature values considered near the decision boundary. By incorporating this range into the decision tree model, more accurate predictions are achieved, resulting in improved illness detection performance.

Based on experiments and analysis, MDVP:Fo(Hz) is identified as the most important feature for the Parkinson dataset, while Glucose and BMI are crucial for the Diabetes dataset. By loosening the left-right boundary for MDVP:Fo(Hz) by 10%, the mistake rate can be reduced to 0.016% without any false positives in Parkinson's disease detection. Similarly, by relaxing the boundaries for BMI and Glucose by 10% simultaneously, the mistake rate for diabetes detection can

be reduced to 0%. These optimization measures enhance the reliability and effectiveness of the illness detection system.

IX. DISCUSSION

In this section, some interesting findings from the experiments will be discussed. Some cases and samples mentioned, like the similarity comparison are made into tables in Appendices.

A. Dataset Limitations and Stability analysis

Both datasets used in this research have notable limitations. The Parkinson dataset suffers from small size and lack of variations, which affects the interpretability of ALE plots. Due to the cancellation of effects in similar samples, ALE tends to display a flat line. On the other hand, PDP focuses on the observed range of feature values, resulting in consistent changes in model predictions. It is important to note that although the majority of PDP and ALE plots show values above the 0 horizontal lines, it does not necessarily imply a positive contribution to the predictions. Rather, it reflects the higher proportion of positive predictions in the dataset.

The diabetes dataset, despite having fewer features, exhibits a more complex correlation among them. The larger size of the decision tree model indicates its attempt to capture complex interactions between features. Strong interactions exist, such as the relationship between pregnancies and age, where higher age suggests a higher possibility of increased pregnancies. Moreover, the relationship between features in this dataset may be nonlinear, while PDP and ALE assume a linear relationship, leading to a conflict.

Despite these limitations, valuable conclusions can still be drawn from the dataset using the feature importance-based methods applied in this research. Permutation importance offers stability by directly observing the impact on model performance, while gini importance utilizes the internal calculations and structure of decision tree algorithms. SHAP is considered the most stable method, satisfying properties like local accuracy, consistency, and missingness. They are less affected by the datasets while PDP and ALE results can be influenced by dataset distribution, model sensitivity to features, and feature interactions. However, they provide insights into how predictions change with variations in feature values.

B. Explanantion Similarty

1) *Global Explanations:* Based on the experimental findings, different feature importance methods tend to produce highly similar explanations within the same machine learning model and dataset, despite employing different calculation theories. This indicates overall consistency among the methods. However, the degree of similarity varies across different ML methods. Random Forest, which combines feature importances from individual trees, generates the most similar explanations to a single Decision Tree. On the other hand, KNN, MLP, and SVM show relatively lower similarity with Decision Trees, while K-means-based classifiers exhibit the lowest similarity. This discrepancy arises due to K-means' insensitivity to small

differences when calculating distances between data points for cluster assignments. When the dataset contains highly similar data points with minimal differences, K-means may struggle to effectively differentiate them, leading to less consistent feature importance rankings (particularly for the Parkinson dataset). Additionally, K-means treats all features equally, disregarding their actual importance or relevance, which can be problematic when the dataset includes strong interactions between features (as in the case of the diabetes dataset).

Drawing meaningful conclusions from PDP and ALE in this situation is challenging. Generally, they lack consistency within a model and compared to other methods. Additionally, the limitations of the adapted datasets hinder a comprehensive understanding of their properties. However, in rare instances when ALE and PDP results align, it indicates the influence of the corresponding feature on the prediction. Furthermore, if multiple ML models consistently emphasize the importance of a particular feature, it can be deemed significant and in line with the feature importance findings.

2) *Local Explanantion:* As observed in the experiments, the similarity of explanations for similar instances within a model is highly dependent on the model's stability, as indicated by the local explanations provided by SHAP. Among the tested models, SVM, KNN, and Random Forest exhibit the most consistent explanations when the instances are similar. On the other hand, Decision Tree shows less stability, especially when a feature value is close to the boundary of a branch, increasing the risk of misclassification. The same conclusion holds true for Decision Tree as an explanation method itself. MLP and K-means-based classifiers are considered the least stable, aligning with the results and conclusions of my experiments.

C. K-means-based Classifier

The implemented K-means-based classifier demonstrates good accuracy with the dataset, achieving 87% accuracy for Parkinson's and 72% accuracy for diabetes. However, this high performance can be attributed, in part, to the imbalanced nature of the dataset. In the case of Parkinson's, where the positive rate in the test set is 82% compared to the overall positive rate of 75%, six out of the seven labeled clusters are assigned a positive label. Similarly, for the Diabetes dataset, with approximately 65% negative rate for both the whole sample and test set, six out of the nine clusters are assigned a negative label. Thus, when the dataset exhibits a high proportion of either the positive or negative class, there is a higher likelihood of clusters being predominantly assigned to that class. Imbalanced classes tend to inflate accuracy. While the classifier demonstrates some predictive power, its performance should be further evaluated on datasets with different distributions and more balanced class data to establish its reliability.

D. Characteristics of Different ML Algorithms

The analysis of feature importance and PDP/ALE results sheds light on how different classifiers assign importance scores. Decision trees prioritize features present in the decision

tree structure, while random forests distribute importance across a wider range of features due to their ensemble nature. The importance distribution of KNN varies based on the similar points it relies on. MLP and SVM show a more uniform decrease in feature importance rankings. On the other hand, the k-means-based classifier assigns relatively high-importance values to all features, treating them equally.

E. Optimization

Multiple experiments were conducted to determine the most influential feature in optimizing decision tree boundaries when implementing the proposed optimization technique. The results indicate that intermediate nodes with higher Gini values play a critical role in the optimization process. Moreover, it is important to consider the key features identified by Gini importance and other explainability methods. When these results align with each other, it often signifies the need to loosen the boundaries of specific nodes. In simpler models, adjusting a single boundary point is typically adequate, but in more complex models, a combination of boundary adjustments may be required.

F. Evaluations for Explainability Methods

- 1) Expressive Power: Decision rules and tree visualization provide transparent and easily understood insights, but may struggle with complex relationships. PDP and ALE capture non-linear and interactive effects, but their explanations may be less intuitive. Feature importance-based methods (SHAP, Gini Importance, Permutation Importance) clearly identify influential features, but may lack rich detail in explanations.
- 2) Translucency: Consistent with the evaluation criteria, my observation aligns with the notion that model-specific methods tend to have higher translucency compared to model-agnostic methods.
- 3) Portability: In line with the evaluation criteria, my observation supports the idea that model-specific methods generally exhibit lower translucency compared to model-agnostic methods.
- 4) Algorithmic Complexity: Most of the methods have a complexity of $O(n \times m)$, where n is the number of instances and m is the number of intervals. The complexity of SHAP depends on the underlying model. For certain models, such as linear models, the complexity can be linear or quadratic in the number of features.

G. Evaluation for explanations

- 1) Fidelity: Based on the Proxy Model Comparison theory and the similarity comparison discussed in the discussion section, Tree-related explanations and explanations for random forests demonstrate the highest fidelity, followed by SVM, MLP, and KNN. K-means-based explanations have the lowest fidelity. Among the explainability methods, tree-related methods show the highest fidelity, while PDP and ALE have lower fidelity compared to permutation importance and SHAP.

- 2) Consistency: The evaluation of consistency is applicable only to model-agnostic methods. As discussed, PDP and ALE exhibit inconsistencies in the two test datasets. In contrast, Permutation importance and SHAP demonstrate better consistency.
- 3) Stability: It was mentioned in the discussion section that the stability heavily relies on the stability of the ML model itself, rather than the choice of explainability method. Notably, decision trees itself as an explanation can show instability when handling boundary values.
- 4) Generalizability: This property was discussed in the experiments and results section. PDP and ALE show limited generalizability, as their performance heavily depends on the dataset characteristics. Similarly, tree-based explanations struggle with generalizability, generating complex and large trees that are challenging to interpret. On the other hand, feature importance-based methods, such as Gini, SHAP, and Permutation, exhibit better generalizability.
- 5) Comprehensibility: ALE and PDP curves are easily understood and indicate the impact of each feature on predictions. Feature importance based measures provide straightforward and easily understandable explanations. Decision rules, tree visualization, and decision rules offer high interpretability, providing transparent insights into how the model makes decisions.
- 6) Degree of Importance: PDP shows the impact of feature changes on model predictions, while ALE illustrates the average effect of each feature. Feature Importance methods measure importance by shuffling feature values. Decision Rules and Tree Visualization offer insights into feature importance.
- 7) Scope: Among the considered methods, Tree-Based explanation and SHAP stand out as they provide both global and local explanations. Other methods primarily offer global explanations without the same level of local interpretability.

X. CONCLUSION

This research delved into the explicitness of various explanation methods employed in different machine learning models. Through a comprehensive evaluation of these methods, we have been able to answer the research question effectively.

- 1) How to measure the utility of the explanation generated by different methods? Evaluation criteria for both the explainability methods and the generated explanations are discussed in Section VI and utilized in Section IX.
 - 2) How do the explainability methods perform given different datasets and what could be the reason if the performance of the methods differs? The results are presented in Section VIII, and further discussed in Section IX.
 - 3) Whether model-specific methods achieve better results than model-agnostic methods? Model-specific methods provide more insightful explanations, while model-agnostic methods offer higher portability. These findings are presented in Sections V, VIII and IX.
- 4) Whether model-agnostic methods could be used on the results of unsupervised learning algorithms? Yes, the introduction of the unsupervised learning algorithm can be found in Section IV, and its implementation is discussed in Section VII and discussion in Section IX.
 - 5) How to improve the performance of the methods? Strategies for improving the performance of the methods are discussed in Section VIII and IX.

A. Future Work

Due to the limitations of the test datasets, the properties of ALE and PDP were not fully examined. Future work could involve running PDP and ALE on more appropriate datasets. Additionally, collecting more practical data with a normal distribution, reflecting the actual illness, would allow for evaluating the suitability of the k-means-based classifier and the proposed decision tree optimization for clinical practice.

REFERENCES

- [1] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, *Machine learning interpretability: A survey on methods and metrics*, *Electronics*, vol. 8, no. 8, pp. 832, 2019.
- [2] M. A. Ahmad, C. Eckert, and A. Teredesai, *Interpretable Machine Learning in Healthcare*, *Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, pp. 559–560, Association for Computing Machinery, 2018.
- [3] A. F. Markus, J. A. Kors, and P. R. Rijnbeek, *The role of explainability in creating trustworthy artificial intelligence for health care: a comprehensive survey of the terminology, design choices, and evaluation strategies*, *Journal of Biomedical Informatics*, vol. 113, pp. 103655, 2021.
- [4] S. R. Islam, W. Eberle, and S. K. Ghafoor, *Towards quantification of explainability in explainable artificial intelligence methods*, *arXiv preprint arXiv:1911.10104*, 2019.
- [5] Max Little, *Parkinsons*, 2008. UCI Machine Learning Repository.
- [6] Smith, J., & Jones, A., *Pima Indians Diabetes*, 1988. UCI Machine Learning Repository.
- [7] Mahesh, Batta. *Machine learning algorithms - a review*. *International Journal of Science and Research (IJSR)*. [Internet], vol. 9, pp. 381–386, 2020.
- [8] Bair, Eric. *Semi-supervised clustering methods*. *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 5, no. 5, pp. 349–361, 2013. Publisher: Wiley Online Library.
- [9] Dasgupta, Sanjoy, Nave Frost, Michal Moshkovitz, and Cyrus Rashtchian. *Explainable k-means clustering: Theory and practice*. In *XXAI Workshop, ICML*, 2020.
- [10] Molnar, Christoph, Timo Freiesleben, Gunnar König, Giuseppe Casalicchio, Marvin N Wright, and Bernd Bischl. *Relating the partial dependence plot and permutation feature importance to the data generating process*. *arXiv preprint arXiv:2109.01433*, 2021.
- [11] Mangalathu, Sujith, Karthika Karthikeyan, De-Cheng Feng, and Jong-Su Jeon. *Machine-learning interpretability techniques for seismic performance assessment of infrastructure systems*. *Engineering Structures*, vol. 250, pp. 112883, 2022. ISSN 0141-0296.
- [12] Selbach, A. (2020). Building an Understandable Explanation System for Machine Learning Models in the Production Domain. Master Thesis DKE-20-21, Maastricht University
- [13] Pedregosa, F. et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- [14] DanaJomar. (2023). *PyALE: Python Accumulated Local Effects*. Retrieved from <https://pypi.org/project/PyALE/>.
- [15] Lundberg, S. (2023). *SHAP: SHapley Additive exPlanations*. Retrieved from <https://pypi.org/project/shap/>.
- [16] Jiangchun, L. (2023). *PDPbox: Partial Dependence Plot Toolbox*. Retrieved from <https://github.com/SauceCat/PDPbox>.

XI. APPENDICES

A. Experimental Settings

Explain Techniques Model	Model-Specific	Model-Agnostic
Decision Tree	Gini Importance, Decision Rules, Tree Visualization	PDP, Permutation Importance, ALE, SHAP (Global & Individual)
Random Forest	Gini Importance	
KNN		
MLP		
SVM		
K-means based model		

Fig. 7. Model Explanation Techniques Table

k (num of clusters)	Accuracy
2	62.99%
3	68.18%
4	68.83%
5	67.53%
6	67.53%
7	70.58%
8	71.42%
9	72.08%
10	71.42%

Fig. 9. Diabetes Choice of Clusters Number for KM-based model

k (num of clusters)	Accuracy
2	82.05%
3	82.05%
4	82.05%
5	84.61%
6	84.61%
7	87.18%
8	74.36%
9	71.79%

Fig. 8. Parkinson Choice of Clusters Number for KM-based model

B. ML Model Performance

Model	Accuracy (%)
DT	92.31
RF	94.87
KNN	82.05
MLP	84.62
SVM	88.14
KM	87.18

(a) Table 1: Parkinson Dataset

Model	Accuracy (%)
DT	79.22
RF	72.08
KNN	76.62
MLP	72.73
SVM	74.46
KM	72.07

(b) Table 2: Diabetes Dataset

Fig. 10. Accuracy of Models on two datasets

C. PDP and ALE in Parkinsons: Conflicts & Agreement

TABLE I
ALE AND PDP: AGREEMENT AND CONFLICTS IN PARKINSON EXPLANATIONS

Feature	DT	RF	KNN	MLP	SVM	KM
MDVP:Fo(Hz)	✓	✗	✓	✗	✗	✗
MDVP:Fhi(Hz)	✓	✗	✗	✓	✗	✗
MDVP:Flo(Hz)	✗	✗	✓	✗	✗	✗
MDVP:Jitter(%)	✓	✗	✓	✗	✗	✗
MDVP:Jitter(Abs)	✓	✗	✓	✗	✗	✗
MDVP:RAP	✗	✗	✓	✗	✗	✗
MDVP:PPQ	✓	✗	✓	✗	✗	✗
Jitter:DDP	✓	✗	✓	✗	✗	✗
MDVP:Shimmer	✓	✗	✓	✗	✗	✗
MDVP:Shimmer(dB)	✓	✗	✓	✗	✗	✗
Shimmer:APQ3	✓	✗	✓	✗	✗	✗
Shimmer:APQ5	✓	✗	✓	✗	✗	✗
MDVP:APQ	✓	✗	✓	✗	✗	✗
Shimmer:DDA	✓	✗	✓	✗	✗	✗
NHR	✓	✗	✓	✗	✗	✗
HNR	✗	✗	✓	✗	✗	✗
RPDE	✗	✗	✗	✗	✗	✗
DFA	✓	✗	✓	✗	✗	✗
spread1	✗	✓	✓	✗	✗	✗
spread2	✓	✗	✓	✗	✗	✗
D2	✗	✗	✓	✗	✗	✗
PPE	✓	✗	✓	✗	✗	✗

'✗' represents both PDP and ALE results conflicts.

'✓' indicates that both PDP and ALE results agree with each other.

TABLE III
CONFLICT BETWEEN ALE AND PDP IN PARKINSONS: TWO SITUATIONS

Feature	DT	RF	KNN	MLP	SVM	KM
MDVP:Fo(Hz)		=		=	=	=
MDVP:Fhi(Hz)		=	✗		=	=
MDVP:Flo(Hz)	=	=		✗	=	=
MDVP:Jitter(%)		=		=	=	=
MDVP:Jitter(Abs)		=		=	=	=
MDVP:RAP	=	=		=	=	=
MDVP:PPQ		=		=	=	=
Jitter:DDP		=		=	=	✗
MDVP:Shimmer		=		=	=	✗
MDVP:Shimmer(dB)		=		=	=	✗
Shimmer:APQ3		=		=	=	✗
Shimmer:APQ5		=		=	=	✗
MDVP:APQ		=		=	=	✗
Shimmer:DDA		=		=	=	=
NHR		=		=	=	=
HNR	=	=		=	=	=
RPDE		=		=	=	=
DFA		=		=	=	=
spread1	=			=	=	=
spread2		=		=	=	=
D2	=	=		=	=	✗
PPE		=		=	=	=

'✗' represents the variation in PDP and ALE results in different directions.

'=' indicates that ALE remains constant at ALE = 0, while PDP exhibits variability.

D. PDP and ALE in Diabetes: Conflicts & Agreement

TABLE IV
ALE AND PDP: AGREEMENT AND CONFLICTS IN DIABETES EXPLANATIONS

Feature	DT	RF	KNN	MLP	SVM	KM
Pregnancies	✗	✗	✗	✗	✗	✓
Glucose	✗	✓	✓	✓	✓	✗
BloodPressure	✗	✗	✗	✓	✗	✓
SkinThickness	✓	✓	✗	✗	✗	✓
Insulin	✗	✗	✗	✗	✗	✓
BMI	✓	✓	✗	✗	✗	✗
DiabetesPedigreeFunction	✓	✓	✗	✗	✗	✓
Age	✓	✓	✗	✓	✓	✓

'✗' represents both PDP and ALE results conflicts.

'✓' indicates that both PDP and ALE results agree with each other.

TABLE II
AGREEMENT BETWEEN PDP AND ALE RESULTS: TREND PATTERNS

Feature	DT	RF	KNN	MLP	SVM	KM
MDVP:Fo(Hz)	Δ		Δ			
MDVP:Fhi(Hz)	=			Δ		
MDVP:Flo(Hz)			Δ			
MDVP:Jitter(%)	=		=			
MDVP:Jitter(Abs)	=		=			
MDVP:RAP			=			
MDVP:PPQ			=			
Jitter:DDP	=		=			
MDVP:Shimmer			=			
MDVP:Shimmer(dB)			=			
Shimmer:APQ3			=			
Shimmer:APQ5			=			
MDVP:APQ			=			
Shimmer:DDA			=			
NHR			=			
HNR			=			
RPDE			=			
DFA			=			
spread1		Δ	=			
spread2		=	=			
D2			=			
PPE	Δ		=			

'=' represents both PDP and ALE results remain consistent along the $y = 0$ line.

'Δ' indicates that both PDP and ALE results vary in the same direction.

E. Influential Features for the Two Datasets

	Permutation	SHAP	Gini
DT	MDVP:Fo(Hz), PPE, MDVP:RAP	MDVP:Fo(Hz), PPE	PPE, MDVP:Fo(Hz)
RF	PPE, spread1	PPE, spread1, MDVP:Fo(Hz)	PPE, MDVP:Fo(Hz), spread1
KNN	MDVP:Flo(Hz), MDVP:Fhi(Hz)	MDVP:Fhi(Hz), MDVP:Flo(Hz), MDVP:Fo(Hz)	/
MLP	MDVP:Fo(Hz)	MDVP:Fhi(Hz), spread1, MDVP:Fo(Hz)	/
SVM	MDVP:Fo(Hz)	MDVP:Fo(Hz), spread1, MDVP:Flo(Hz), D2, PPE	/
KM	MDVP:Fo(Hz), spread1	D2, MDVP:Flo(Hz), spread1, MDVP:Fo(Hz), PPE, MDVP:Fhi(Hz)	/
Influential features are defined as those with importance scores greater than the non-zero median feature importance score			

Fig. 11. Parkinson Influential Feature Identified by Different ML Models and Different Explan Methods

	Permutation	SHAP	Gini
DT	Glucose, BMI	Glucose, BMI, Age	Glucose, BMI
RF	Glucose, BMI	Glucose, BMI	Glucose
KNN	Glucose	Glucose, Insulin, Age	/
MLP	Glucose, Age, Insulin	Glucose	/
SVM	Glucose	Glucose, Age, BMI	/
KM	Age, SkinThickness	Glucose, Age, BMI	/
Influential features are defined as those with importance scores greater than the non-zero median feature importance score			

Fig. 12. Diabetes Influential Feature Identified by Different ML Models and Different Explan Methods

F. Global SHAP Explanations for Parkinsons

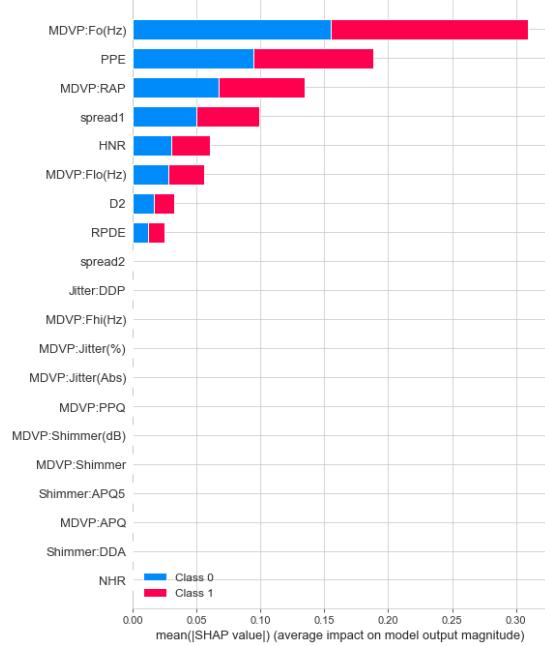


Fig. 13. Global SHAP Decision Tree

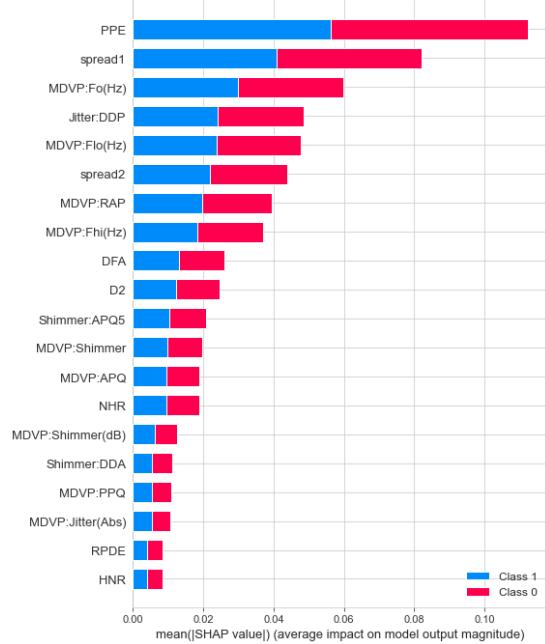


Fig. 14. Global SHAP Decision Tree

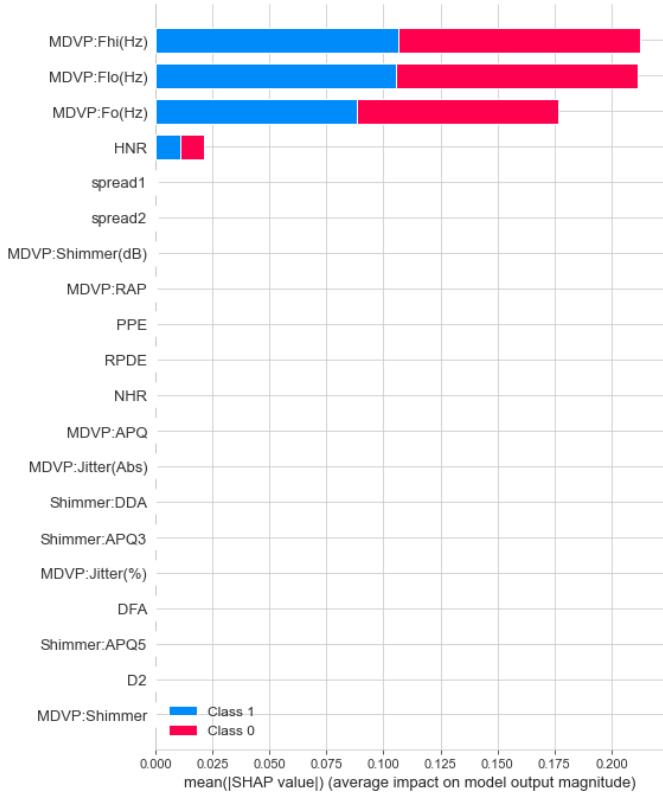


Fig. 15. Global SHAP Random Forest

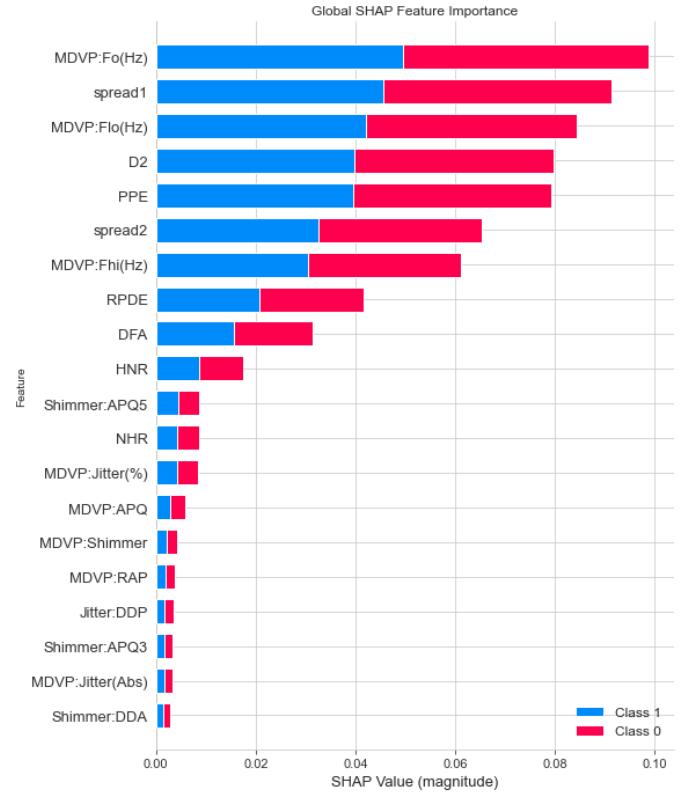


Fig. 17. Global SHAP SVM

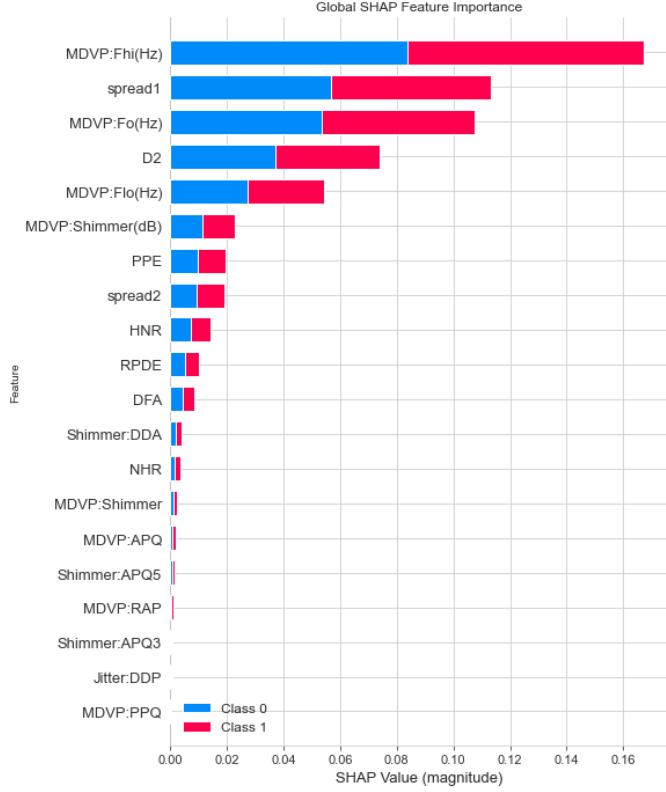


Fig. 16. Global SHAP MLP

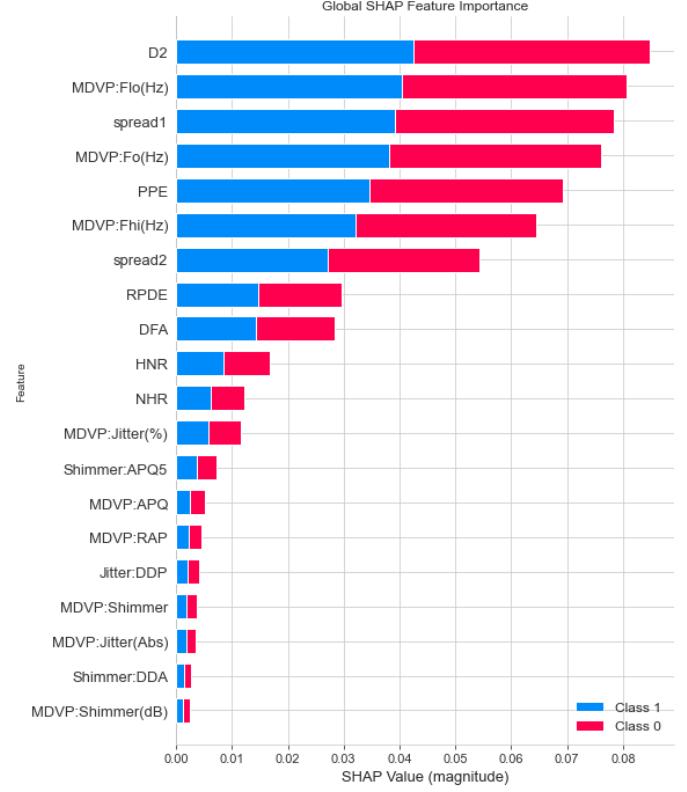


Fig. 18. Global SHAP KM

G. Two Similar instances mentioned in Experiments Parkinsons

For the first data point (phon_R01_S34_4):

MDVP:Fo(Hz): 202.80500
 MDVP:Fhi(Hz): 231.50800
 MDVP:Flo(Hz): 86.23200
 MDVP:Jitter(%): 0.00370
 MDVP:Jitter(Abs): 0.00002
 MDVP:RAP: 0.00189
 MDVP:PPQ: 0.00211
 Jitter:DDP: 0.00568
 MDVP:Shimmer: 0.01997
 MDVP:Shimmer(dB): 0.18000
 Shimmer:APQ3: 0.01117
 Shimmer:APQ5: 0.01177
 MDVP:APQ: 0.01506
 Shimmer:DDA: 0.03350
 NHR: 0.02010
 HNR: 18.68700
 RPDE: 1
 DFA: 0.536102
 spread1: 0.632631
 spread2: -5.898673
 D2: 0.213353
 PPE: 2.470746

Fig. 19. Parkinson 113

For the second data point (phon_R01_S26_5):

MDVP:Fo(Hz): 210.14100
 MDVP:Fhi(Hz): 232.70600
 MDVP:Flo(Hz): 185.25800
 MDVP:Jitter(%): 0.00534
 MDVP:Jitter(Abs): 0.00003
 MDVP:RAP: 0.00321
 MDVP:PPQ: 0.00280
 Jitter:DDP: 0.00964
 MDVP:Shimmer: 0.01680
 MDVP:Shimmer(dB): 0.14900
 Shimmer:APQ3: 0.00861
 Shimmer:APQ5: 0.01017
 MDVP:APQ: 0.01301
 Shimmer:DDA: 0.02583
 NHR: 0.00620
 HNR: 23.67100
 RPDE: 1
 DFA: 0.441097
 spread1: 0.722254
 spread2: -5.963040
 D2: 0.250283
 PPE: 2.489191

Fig. 20. Parkinson 144

H. Global SHAP Explanations for Diabetes

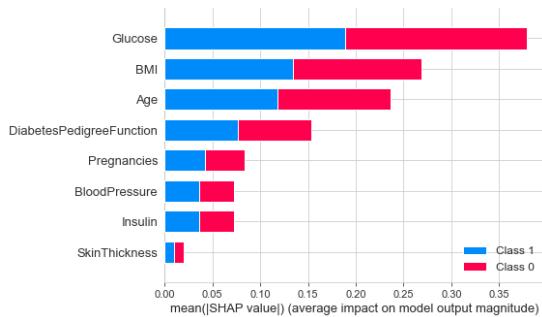


Fig. 21. Global SHAP Decision Tree

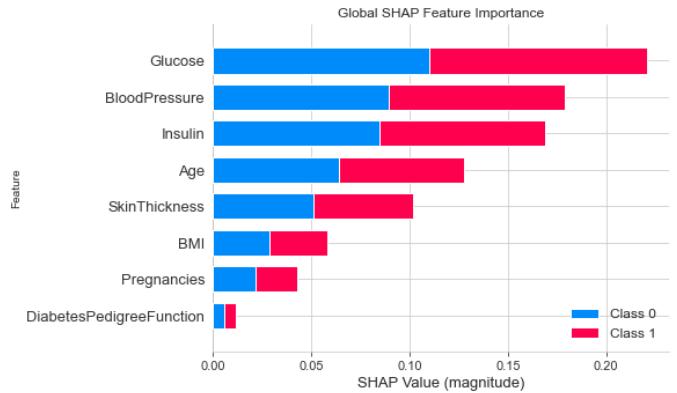


Fig. 24. Global SHAP MLP

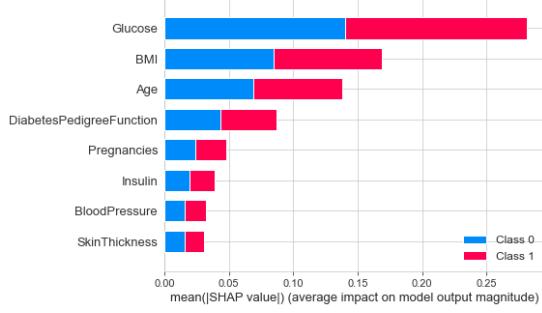


Fig. 22. Global SHAP Decision Tree

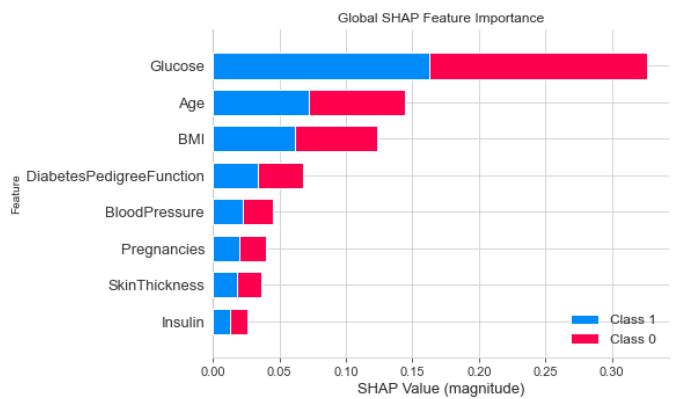


Fig. 25. Global SHAP SVM

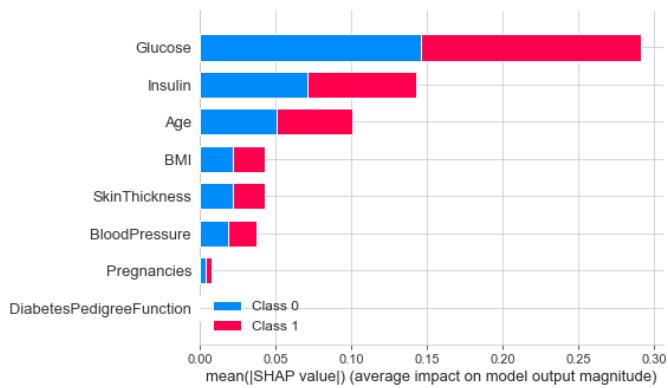


Fig. 23. Global SHAP Random Forest

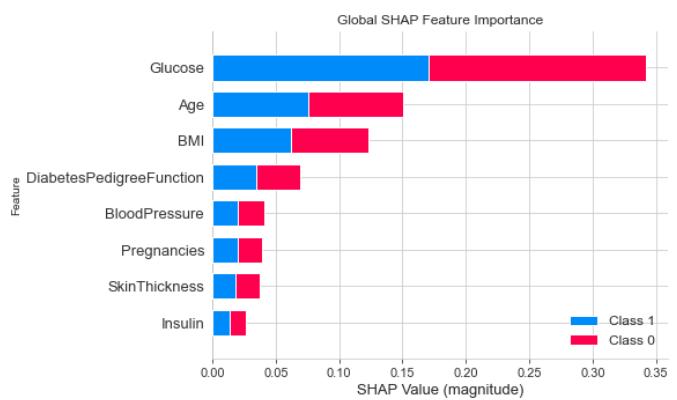


Fig. 26. Global SHAP KM

I. Permutation Importance Explanations for Parkinsons

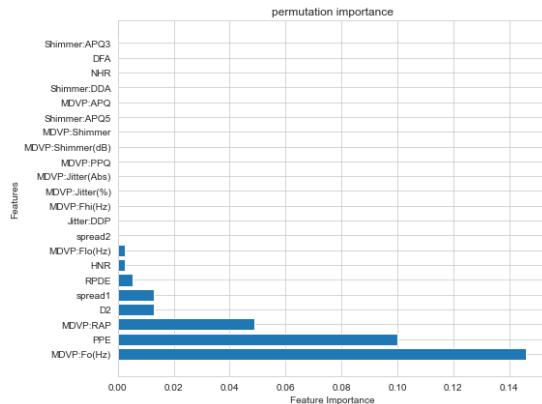


Fig. 27. Permutation Importance Decision Tree

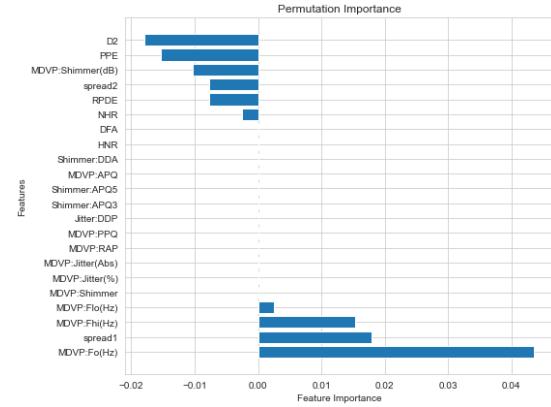


Fig. 30. Permutation Importance MLP

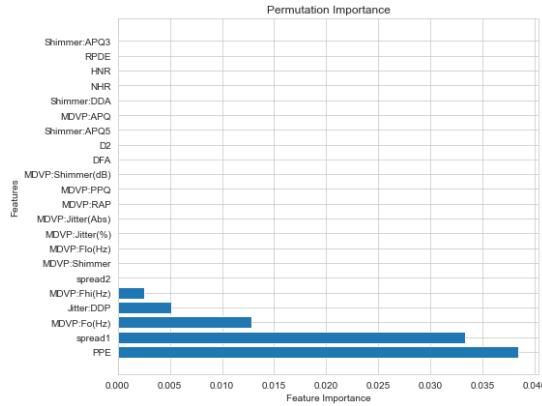


Fig. 28. Permutation Importance Decision Tree

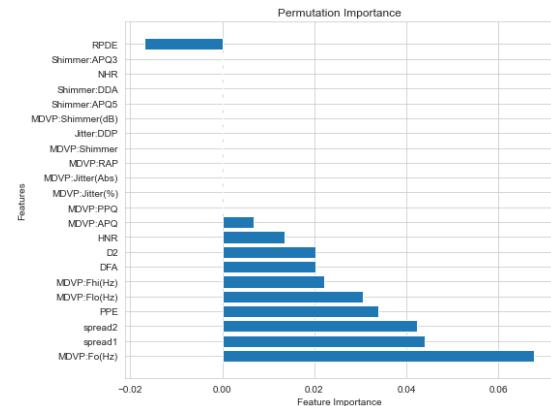


Fig. 31. Permutation Importance SVM

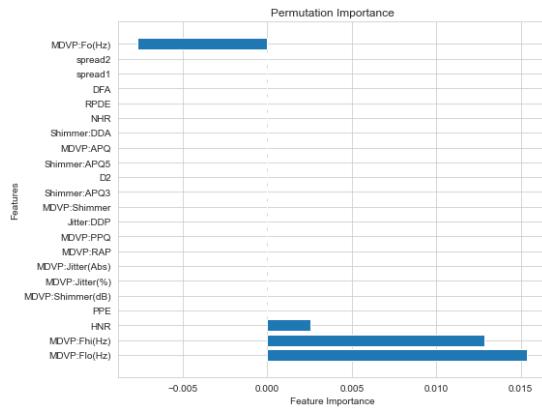


Fig. 29. Permutation Importance Random Forest

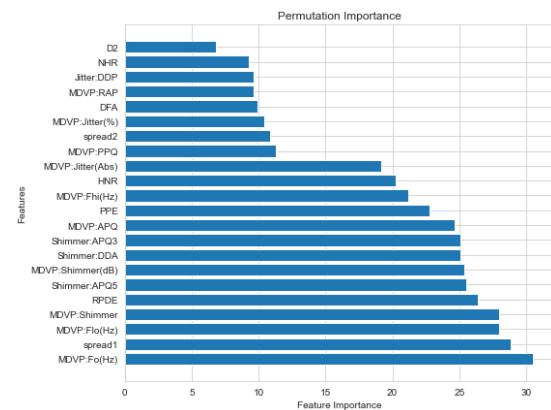


Fig. 32. Permutation Importance KM

J. Permutation Importance Explanations for Parkinsons

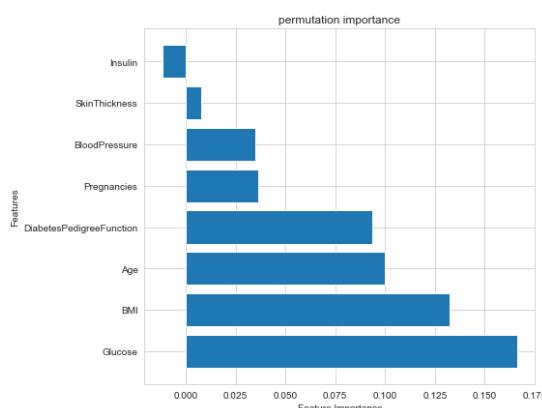


Fig. 33. Permutation Importance Decision Tree

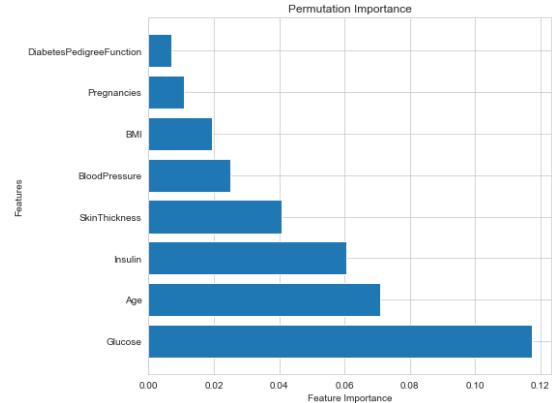


Fig. 36. Permutation Importance MLP

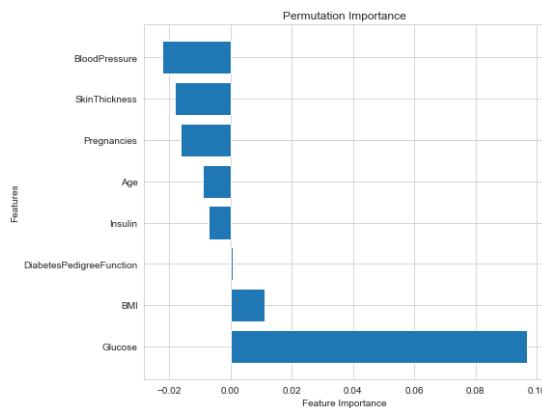


Fig. 34. Permutation Importance Decision Tree

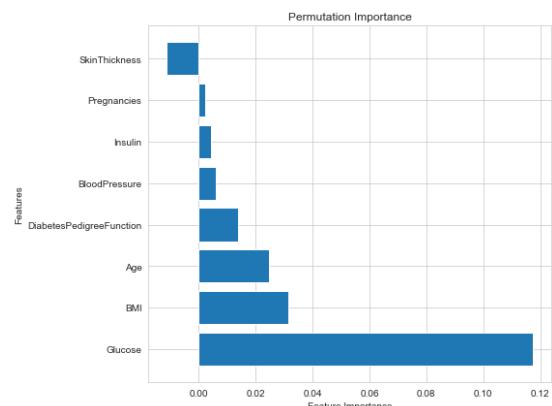


Fig. 37. Permutation Importance SVM

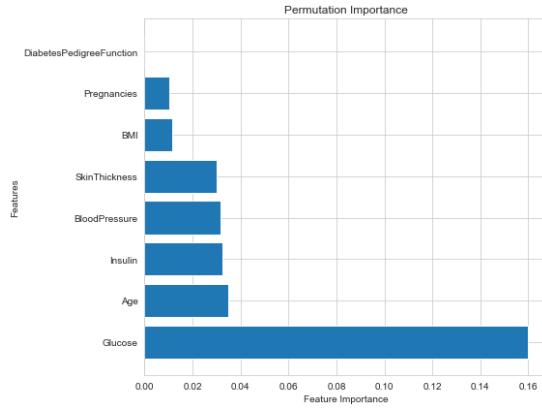


Fig. 35. Permutation Importance Random Forest

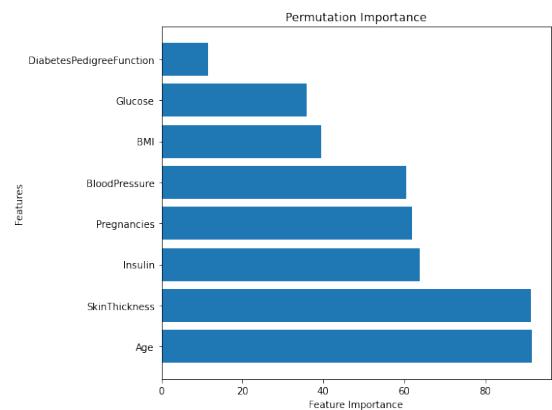


Fig. 38. Permutation Importance KM

K. Tree Based Explanations for Parkinsons

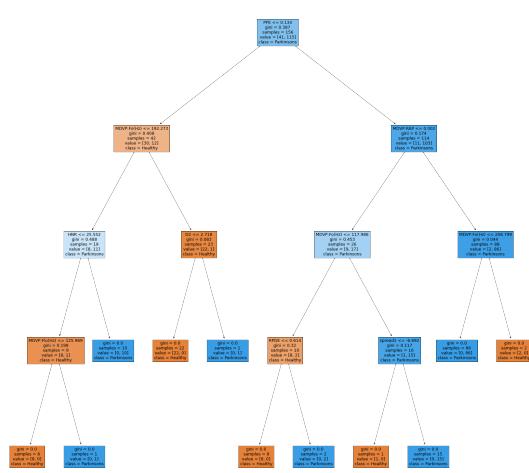


Fig. 39. Decision Tree Visualization

```

--- PPE <= 0.13
|--- MDVP:Fo(Hz) <= 192.27
|   |--- HNR <= 25.55
|   |   |--- MDVP:F1o(Hz) <= 125.97
|   |   |   |--- class: 0
|   |   |   |--- MDVP:F1o(Hz) > 125.97
|   |   |   |--- class: 1
|   |--- HNR > 25.55
|   |   |--- class: 1
|--- MDVP:Fo(Hz) > 192.27
|   |--- D2 <= 2.72
|   |   |--- class: 0
|   |--- D2 > 2.72
|   |   |--- class: 1
--- PPE > 0.13
|--- MDVP:RAP <= 0.00
|   |--- MDVP:Fo(Hz) <= 117.99
|   |   |--- RPDE <= 0.61
|   |   |   |--- class: 0
|   |   |   |--- RPDE > 0.61
|   |   |   |--- class: 1
|--- MDVP:fo(Hz) > 117.99
|   |--- spread1 <= -6.69
|   |   |--- class: 0
|   |   |--- spread1 > -6.69
|   |   |--- class: 1
|--- MDVP:RAP > 0.00
|   |--- MDVP:Fo(Hz) <= 208.80
|   |   |--- class: 1
|   |--- MDVP:fo(Hz) > 208.80
|   |   |--- class: 0

```

Fig. 40. Decision Rule

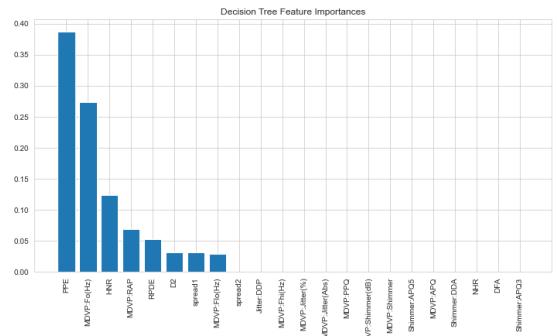


Fig. 41. Gini Importance

L. Tree Based Explanations for Diabetes

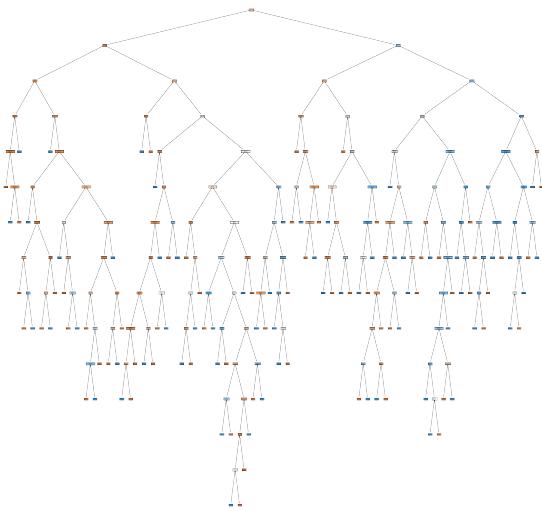


Fig. 42. Decision Tree Visualization

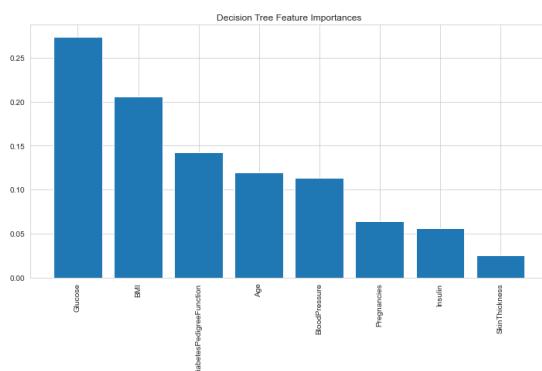


Fig. 43. Gini Importance

```

--- Glucose <= 127.50
|--- Age <= 28.50
|--- BMI <= 30.95
|--- Pregnancies <= 7.50
|   |--- DiabetesPedigreeFunction <= 0.67
|   |--- class: 0
|   |--- DiabetesPedigreeFunction > 0.67
|   |--- class: 1
|   |--- DiabetesPedigreeFunction > 0.69
|   |--- class: 0
|--- Pregnancies > 7.50
|   |--- class: 1
|--- BMI > 30.95
|--- BloodPressure <= 37.00
|   |--- class: 1
|--- BloodPressure > 37.00
|   |--- DiabetesPedigreeFunction <= 0.50
|   |--- BMI <= 31.10
|   |--- class: 1
|   |--- BMI > 31.10
|       |--- SkinThickness <= 8.00
|       |--- Glucose <= 100.50
|       |--- class: 0
|       |--- Glucose > 100.50
|       |--- BMI <= 32.45
|       |--- class: 0
|       |--- BMI > 32.45
|       |--- class: 1
|       |--- SkinThickness > 8.00
|       |--- BMI <= 31.80
|       |--- Age <= 24.50
|       |--- class: 0
|       |--- Age > 24.50
|       |--- class: 1
|       |--- BMI > 31.80
|       |--- class: 0
|--- DiabetesPedigreeFunction > 0.50
|--- DiabetesPedigreeFunction <= 0.54
|--- BMI <= 34.00
|   |--- class: 1
|--- BMI > 34.00
|--- Pregnancies <= 1.50
|   |--- class: 0
|--- Pregnancies > 1.50
|   |--- SkinThickness <= 36.50
|   |--- class: 0
|   |--- SkinThickness > 36.50
|   |--- class: 1
|--- DiabetesPedigreeFunction > 0.54
|--- DiabetesPedigreeFunction <= 1.27
|--- BloodPressure <= 67.00
|--- Insulin <= 70.00
|   |--- class: 0
|--- Insulin > 70.00
|   |--- Glucose <= 114.00
|       |--- DiabetesPedigreeFunction <= 0.64
|       |--- class: 0
|       |--- DiabetesPedigreeFunction > 0.64
|       |--- class: 1
|   |--- Glucose > 114.00
|       |--- class: 0
|--- BloodPressure > 67.00
|   |--- Age <= 21.50
|       |--- BMI <= 42.80
|       |--- class: 0
|       |--- BMI > 42.80
|       |--- class: 1
|   |--- Age > 21.50
|       |--- class: 0
|--- DiabetesPedigreeFunction > 1.27
|   |--- class: 1
|--- Age > 28.50
|--- BMI <= 26.35
|--- BMI <= 0.65
|   |--- class: 1
|--- BMI > 26.35
|   |--- class: 0
|--- BMI > 26.35
|--- Glucose <= 99.50
|--- Glucose <= 28.50
|   |--- class: 1
|--- Glucose > 28.50
|--- Insulin <= 87.00
|   |--- DiabetesPedigreeFunction <= 1.16
|       |--- Age <= 51.00
|       |--- Pregnancies <= 11.50
|       |--- DiabetesPedigreeFunction <= 0.18
|       |--- class: 0
|       |--- DiabetesPedigreeFunction > 0.18
|       |--- class: 1
|   |--- BMI <= 32.40
|       |--- class: 1
|       |--- BMI > 32.40
|       |--- class: 0
|       |--- DiabetesPedigreeFunction > 0.18
|       |--- class: 0
|--- Pregnancies > 11.50
|--- BMI <= 31.25
|   |--- class: 1
|--- BMI > 31.25
|   |--- class: 0
|--- Age > 51.00
|--- SkinThickness <= 13.50
|   |--- class: 0
|--- SkinThickness > 13.50
|   |--- class: 1
|--- DiabetesPedigreeFunction > 1.16
|   |--- class: 1
|--- Insulin > 87.00
|--- BMI <= 33.50
|   |--- class: 0
|--- BMI > 33.50
|   |--- class: 1
|--- Glucose > 99.50
|--- DiabetesPedigreeFunction <= 0.52
|--- DiabetesPedigreeFunction <= 0.20
|   |--- Age <= 34.50

```

Fig. 44. Decision Rule

M. ALE and PDP plots for the Decision Tree on the Parkinson

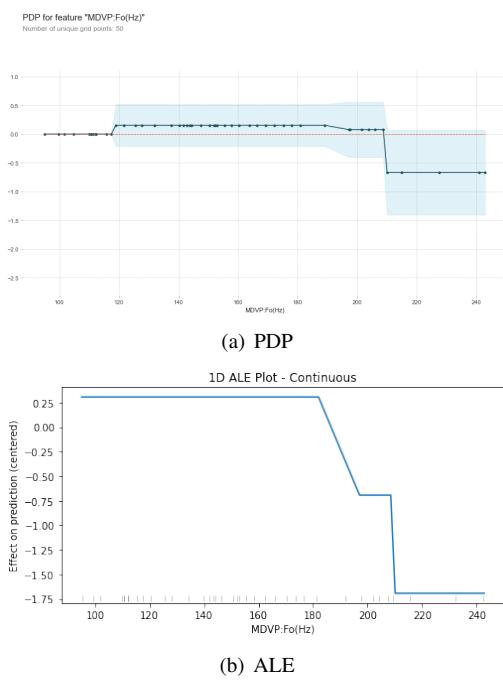


Fig. 45. PDP and ALE for the "MDVP:Fo(Hz)" Feature in the Decision Tree

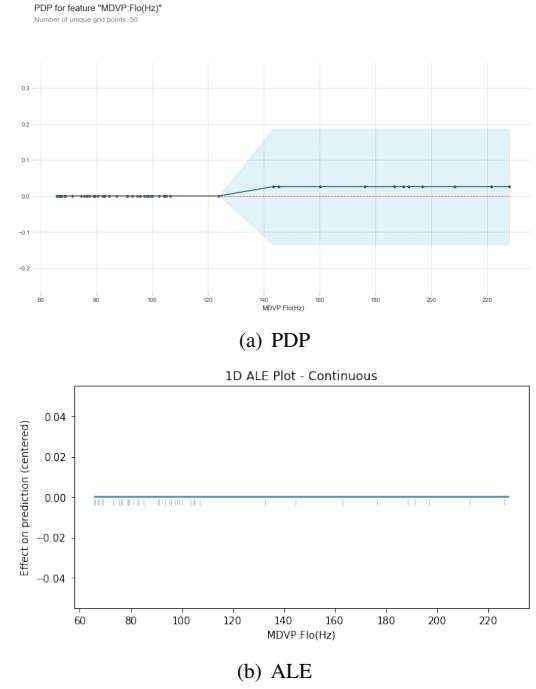


Fig. 47. PDP and ALE for the "MDVP:Flo(Hz)" Feature in the Decision Tree

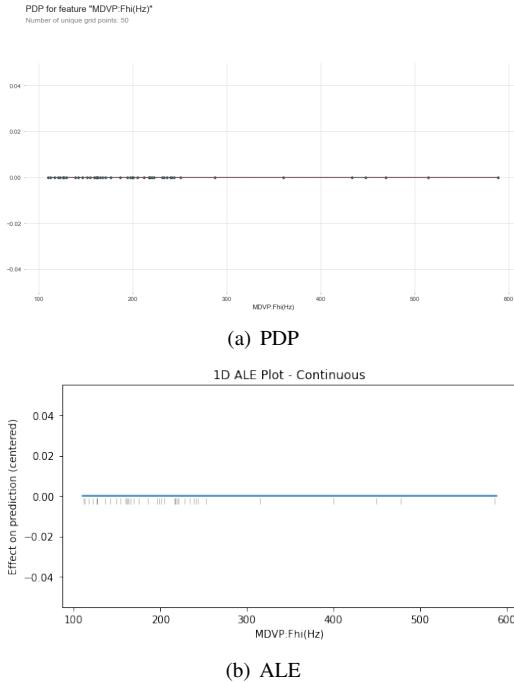


Fig. 46. PDP and ALE for the "MDVP:Fhi(Hz)" Feature in the Decision Tree

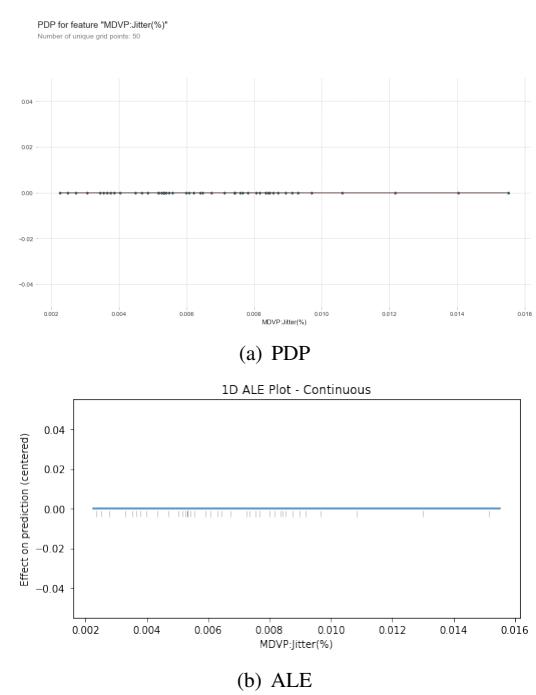
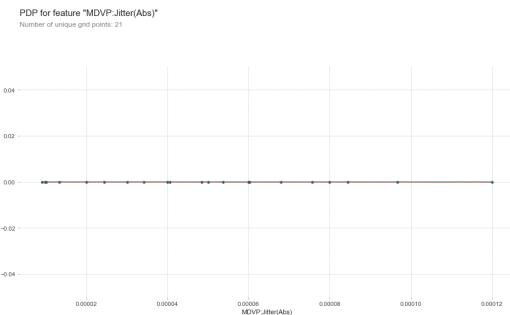
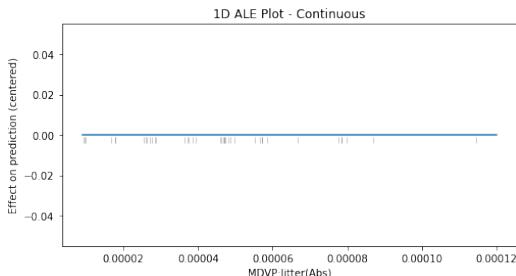


Fig. 48. PDP and ALE for the "MDVP:Jitter(%)" Feature in the Decision Tree

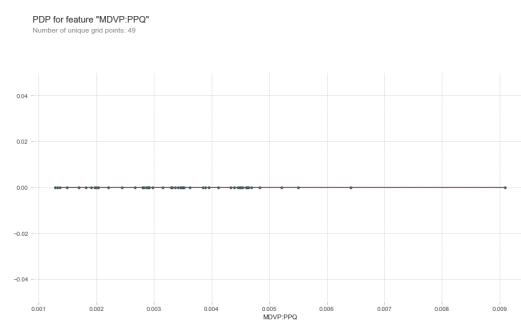


(a) PDP

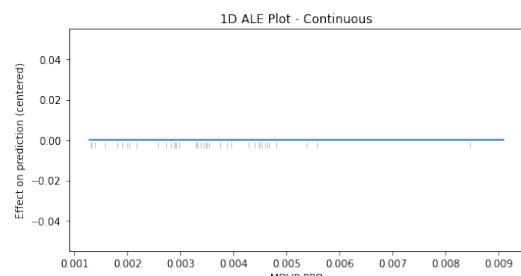


(b) ALE

Fig. 49. PDP and ALE for the "MDVP:Jitter(Abs)" Feature in the Decision Tree

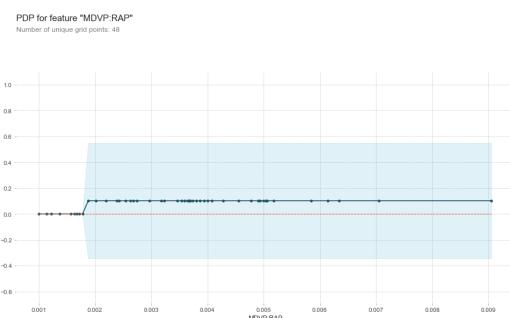


(a) PDP

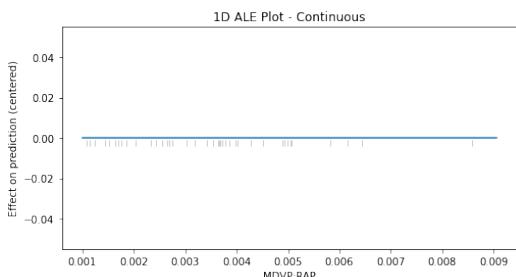


(b) ALE

Fig. 51. PDP and ALE for the "MDVP:PPQ" Feature in the Decision Tree

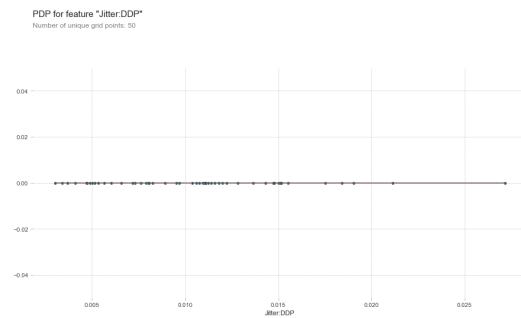


(a) PDP

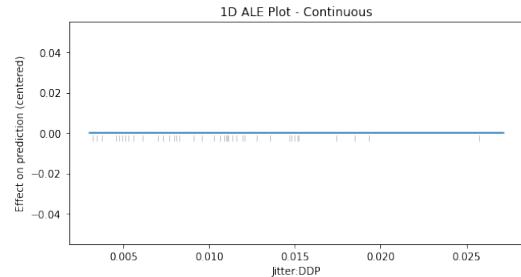


(b) ALE

Fig. 50. PDP and ALE for the "MDVP:RAP" Feature in the Decision Tree



(a) PDP



(b) ALE

Fig. 52. PDP and ALE for the "Jitter:DDP" Feature in the Decision Tree

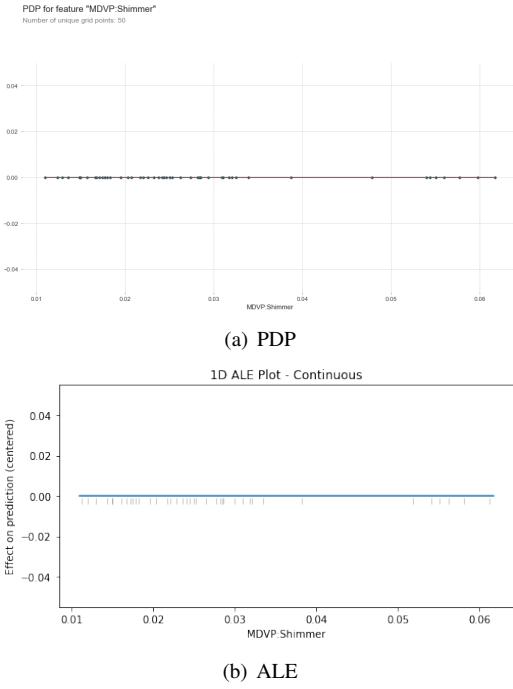


Fig. 53. PDP and ALE for the "MDVP:Shimmer" Feature in the Decision Tree

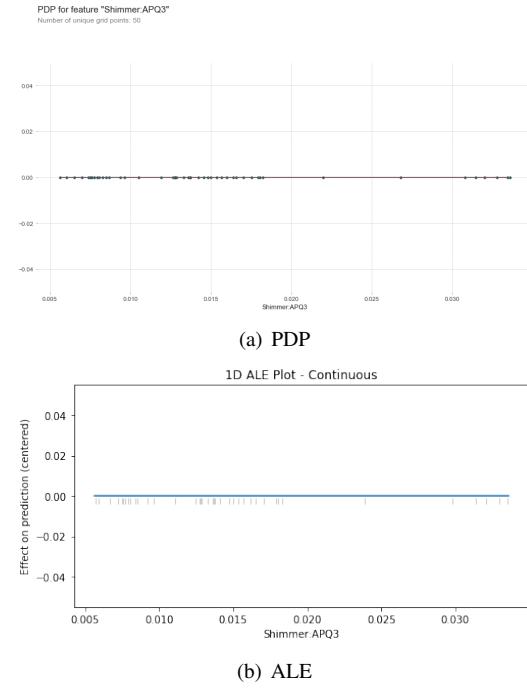


Fig. 55. PDP and ALE for the "Shimmer:APQ3" Feature in the Decision Tree

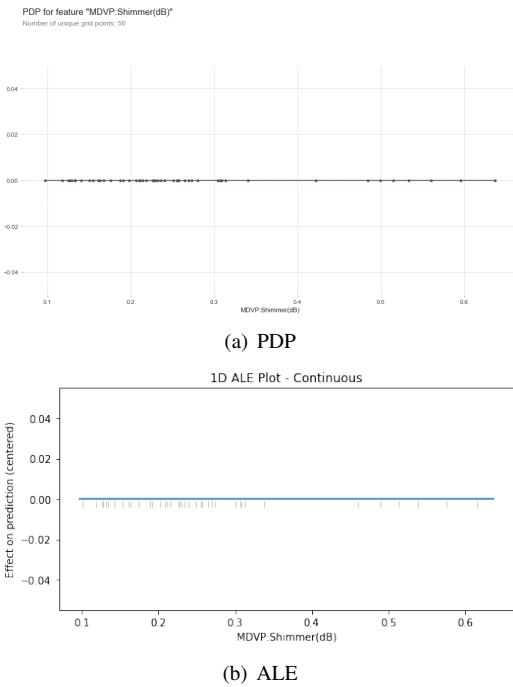


Fig. 54. PDP and ALE for the "MDVP:Shimmer(dB)" Feature in the Decision Tree

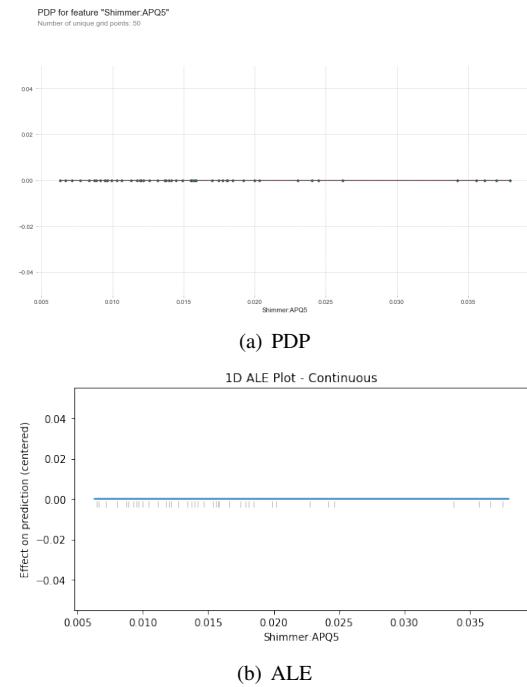
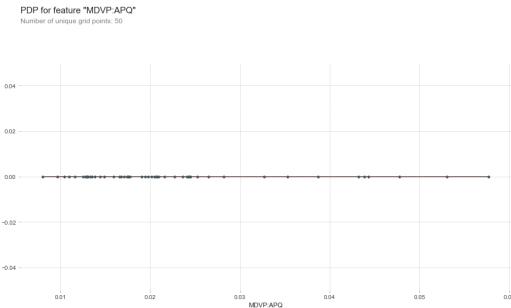
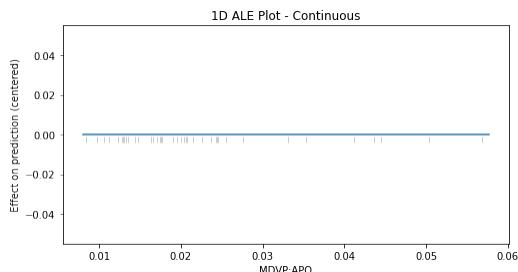


Fig. 56. PDP and ALE for the "Shimmer:APQ5" Feature in the Decision Tree

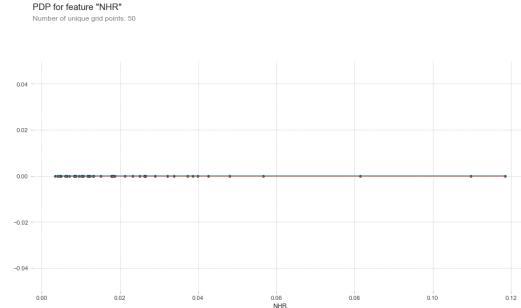


(a) PDP

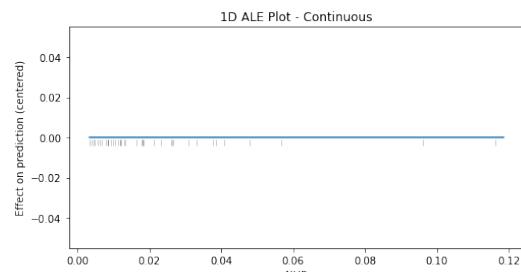


(b) ALE

Fig. 57. PDP and ALE for the "MDVP:APQ" Feature in the Decision Tree

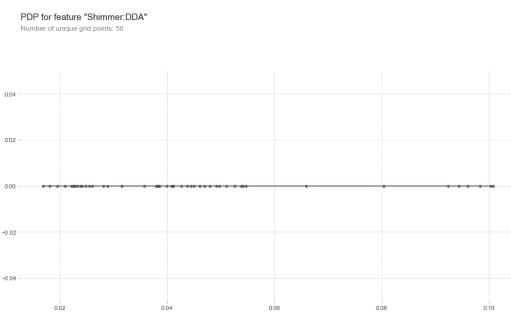


(a) PDP

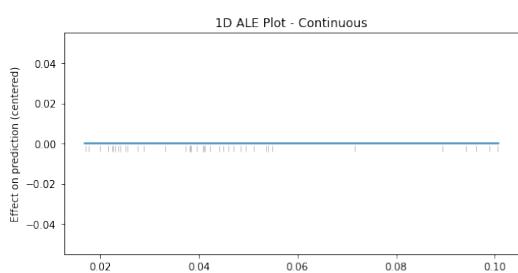


(b) ALE

Fig. 59. PDP and ALE for the "NHR" Feature in the Decision Tree

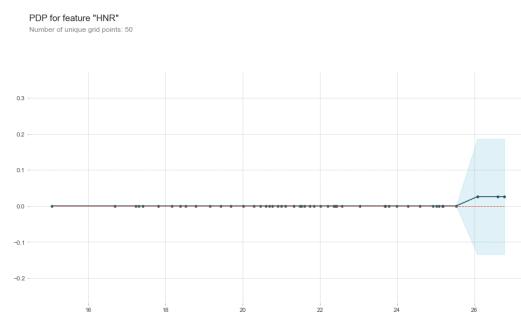


(a) PDP

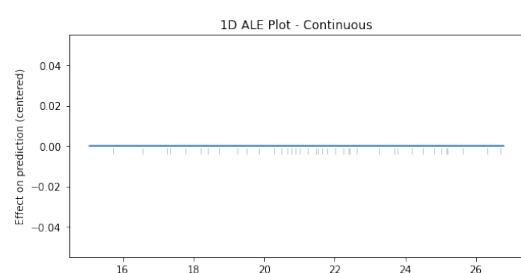


(b) ALE

Fig. 58. PDP and ALE for the "Shimmer:DDA" Feature in the Decision Tree

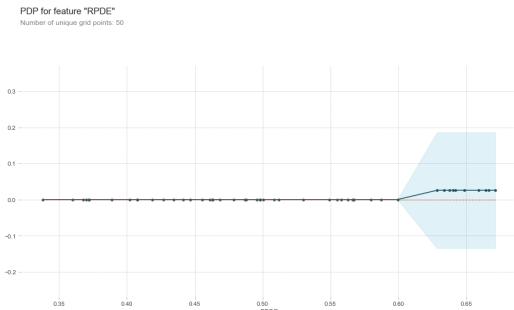


(a) PDP

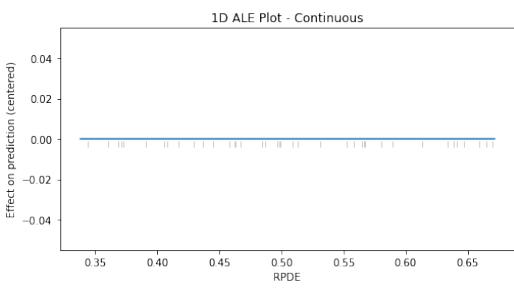


(b) ALE

Fig. 60. PDP and ALE for the "HNR" Feature in the Decision Tree

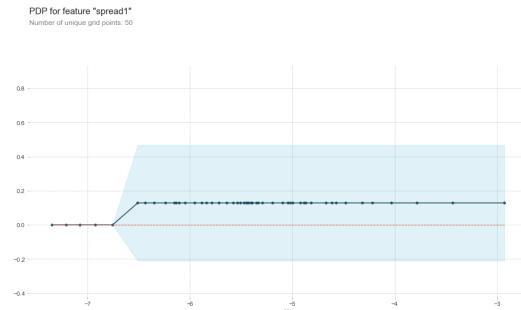


(a) PDP

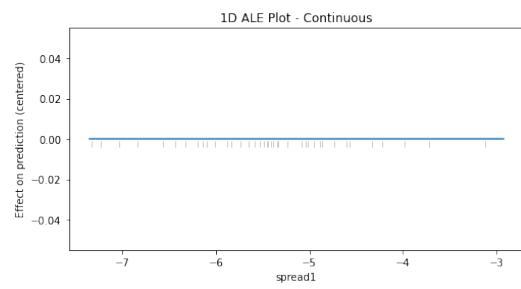


(b) ALE

Fig. 61. PDP and ALE for the "RPDE" Feature in the Decision Tree

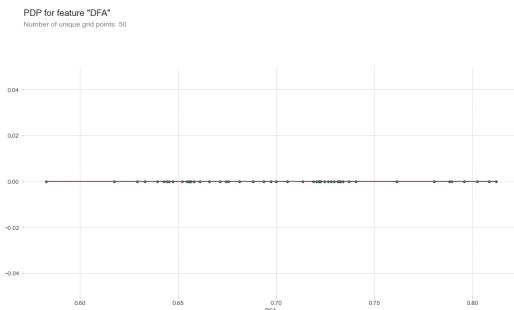


(a) PDP

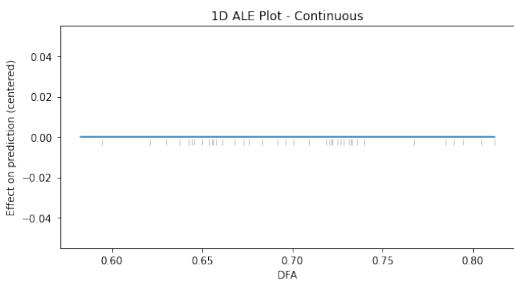


(b) ALE

Fig. 63. PDP and ALE for the "spread1" Feature in the Decision Tree

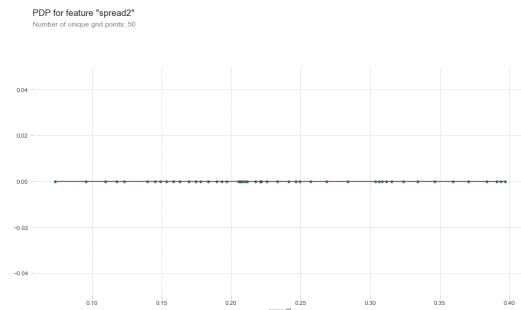


(a) PDP

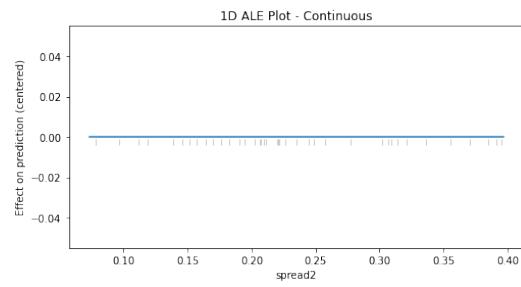


(b) ALE

Fig. 62. PDP and ALE for the "DFA" Feature in the Decision Tree



(a) PDP



(b) ALE

Fig. 64. PDP and ALE for the "spread2" Feature in the Decision Tree

N. ALE and PDP plots for the Random Forest on the Parkinson

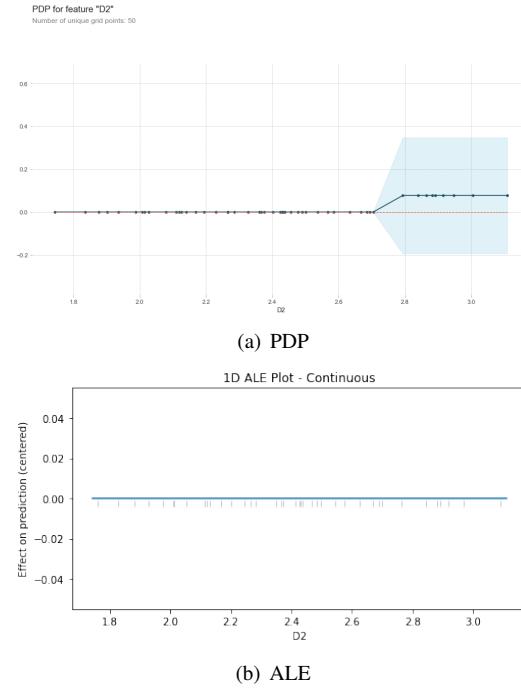


Fig. 65. PDP and ALE for the "D2" Feature in the Decision Tree

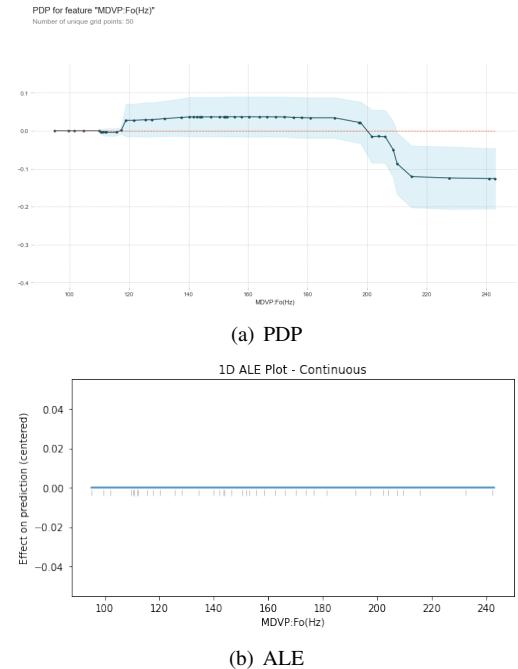


Fig. 67. PDP and ALE for the "MDVP:Fo(Hz)" Feature in the Random Forest

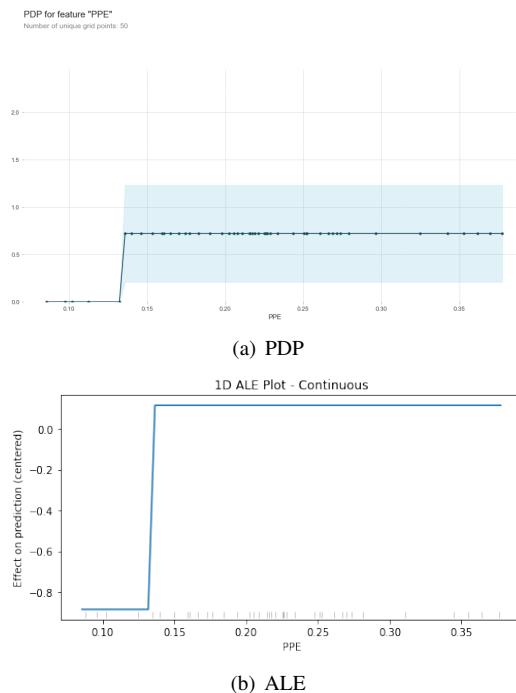


Fig. 66. PDP and ALE for the "PPE" Feature in the Decision Tree

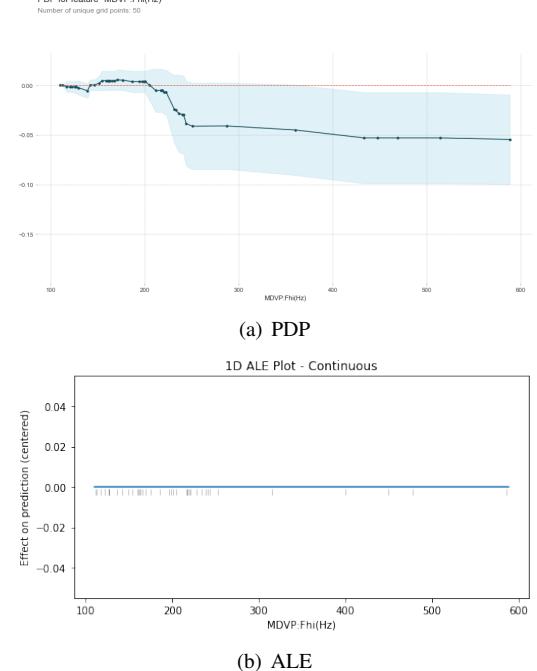
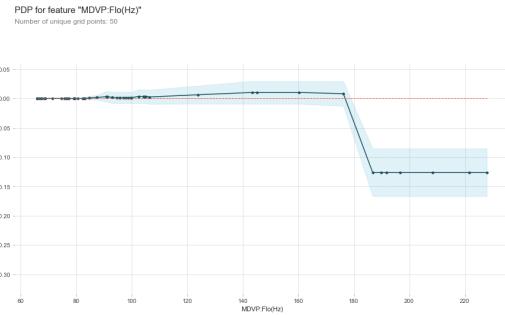
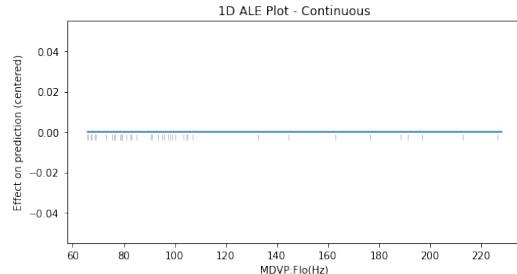


Fig. 68. PDP and ALE for the "MDVP:Fhi(Hz)" Feature in the Random Forest

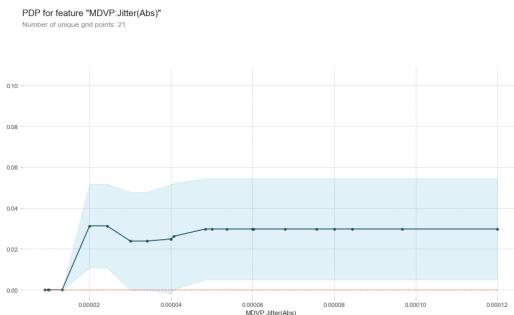


(a) PDP

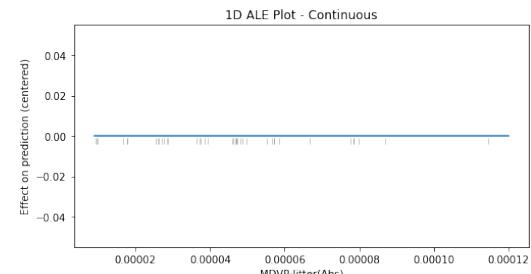


(b) ALE

Fig. 69. PDP and ALE for the "MDVP:Flo(Hz)" Feature in the Random Forest

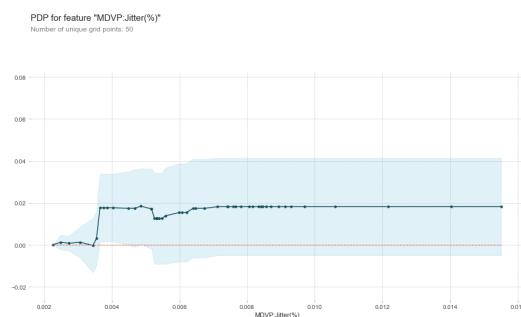


(a) PDP

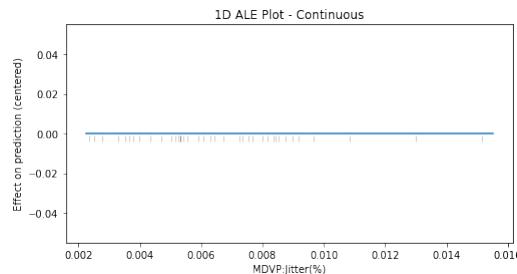


(b) ALE

Fig. 71. PDP and ALE for the "MDVP:Jitter(Abs)" Feature in the Random Forest

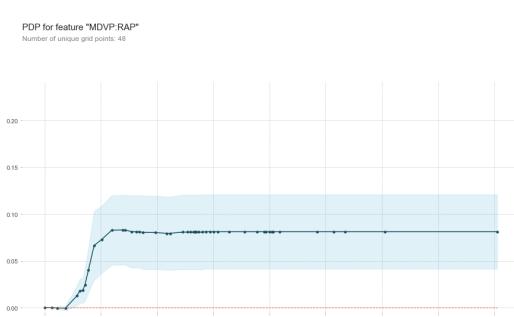


(a) PDP

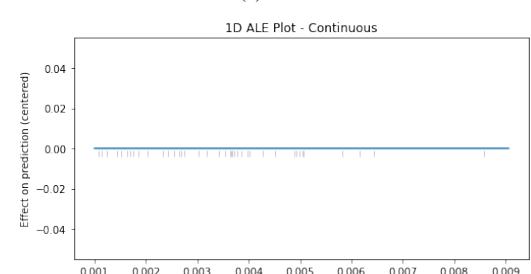


(b) ALE

Fig. 70. PDP and ALE for the "MDVP:Jitter(%)" Feature in the Random Forest

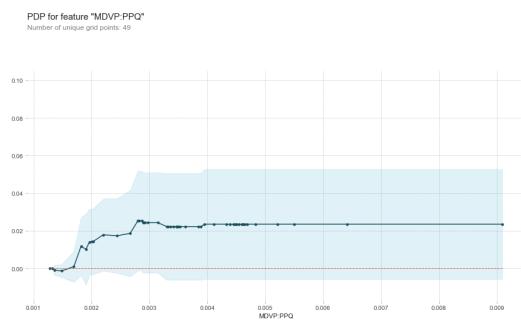


(a) PDP

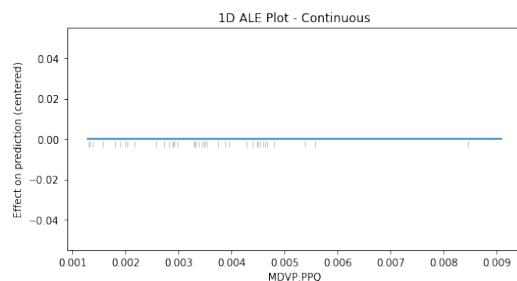


(b) ALE

Fig. 72. PDP and ALE for the "MDVP:RAP" Feature in the Random Forest

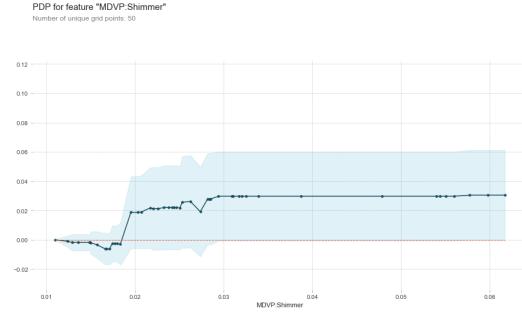


(a) PDP

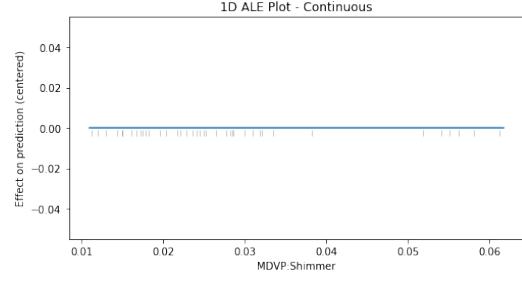


(b) ALE

Fig. 73. PDP and ALE for the "MDVP:PPQ" Feature in the Random Forest

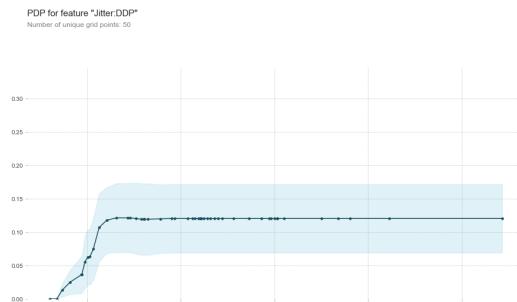


(a) PDP

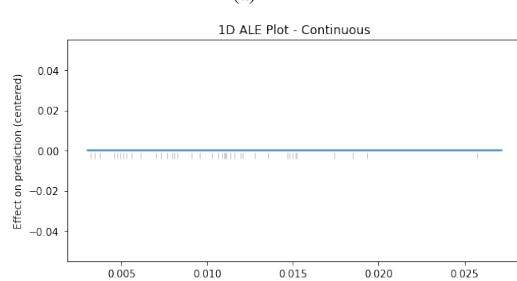


(b) ALE

Fig. 75. PDP and ALE for the "MDVP:Shimmer" Feature in the Random Forest

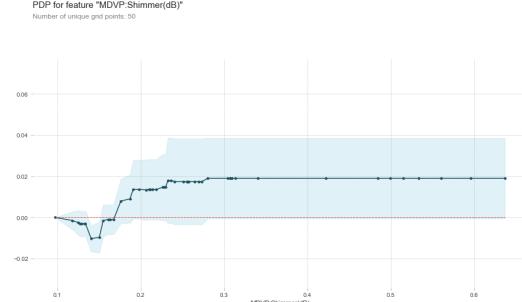


(a) PDP

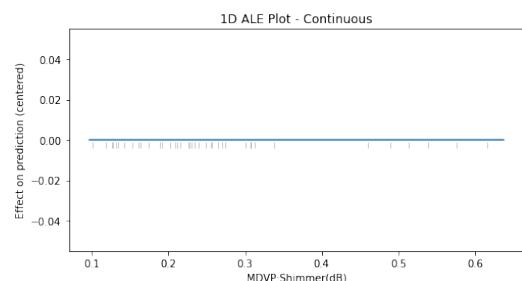


(b) ALE

Fig. 74. PDP and ALE for the "Jitter:DDP" Feature in the Random Forest



(a) PDP



(b) ALE

Fig. 76. PDP and ALE for the "MDVP:Shimmer(dB)" Feature in the Random Forest

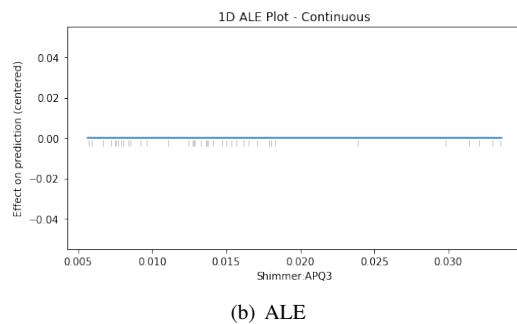
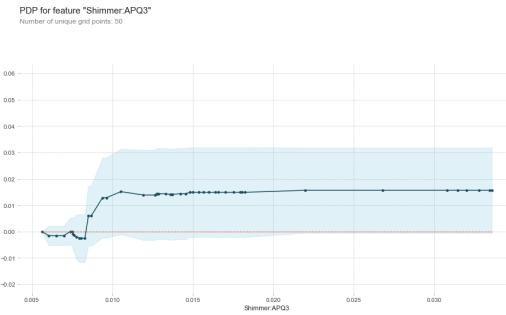


Fig. 77. PDP and ALE for the "Shimmer:APQ3" Feature in the Random Forest

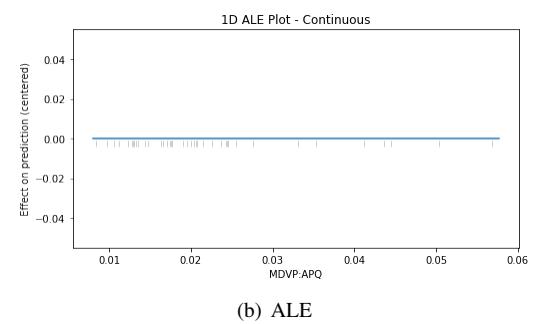
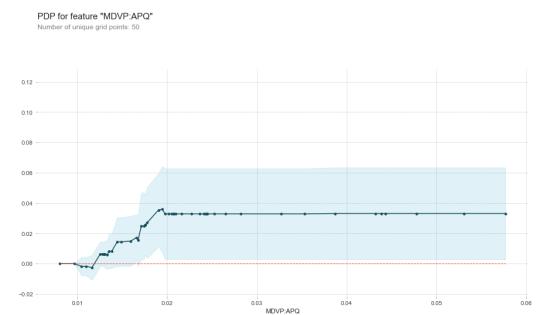


Fig. 79. PDP and ALE for the "MDVP:APQ" Feature in the Random Forest

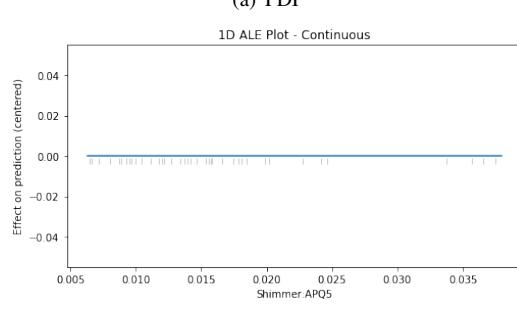
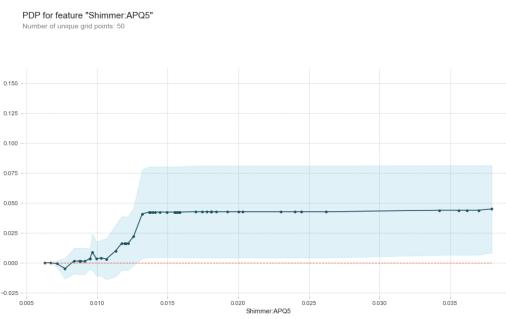


Fig. 78. PDP and ALE for the "Shimmer:APQ5" Feature in the Random Forest

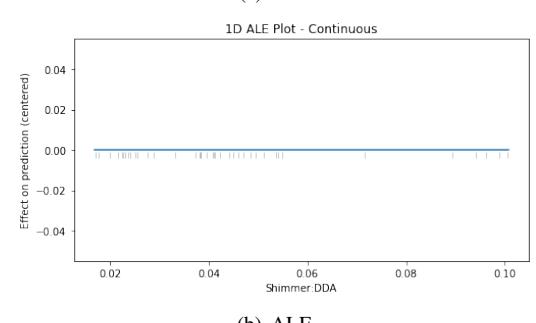
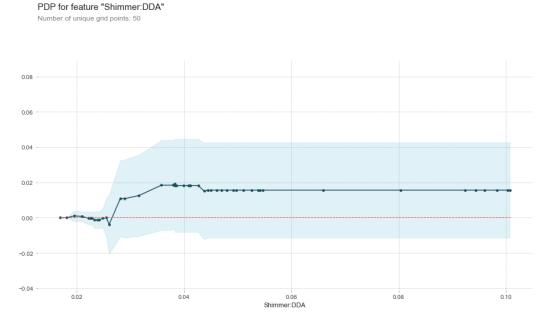


Fig. 80. PDP and ALE for the "Shimmer:DDA" Feature in the Random Forest

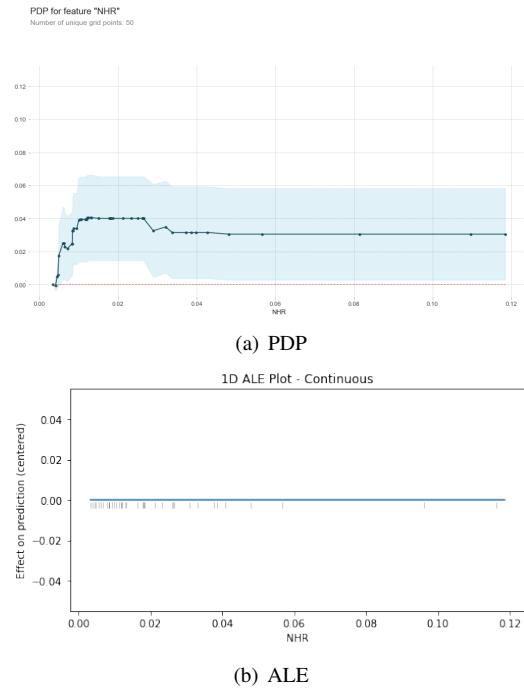


Fig. 81. PDP and ALE for the "NHR" Feature in the Random Forest

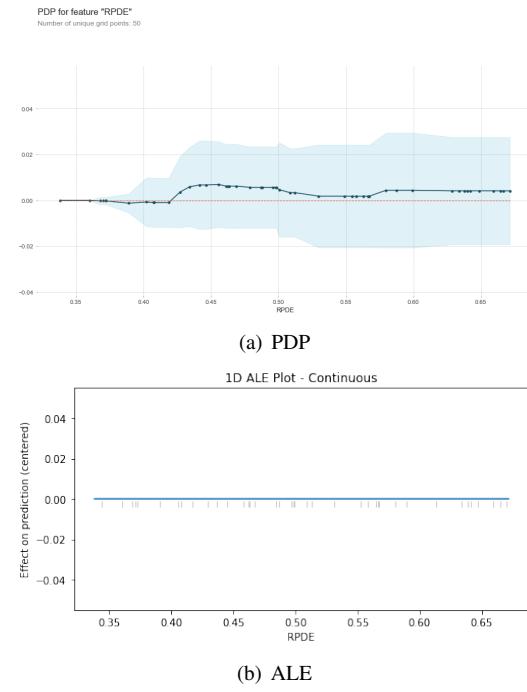


Fig. 83. PDP and ALE for the "RPDE" Feature in the Random Forest

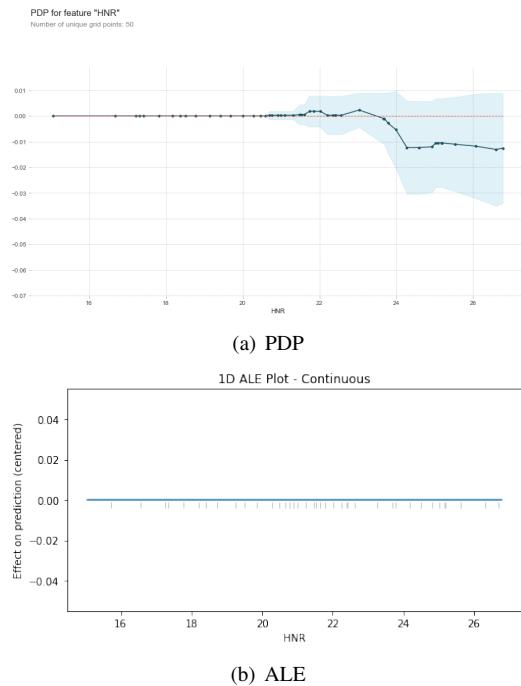


Fig. 82. PDP and ALE for the "HNR" Feature in the Random Forest

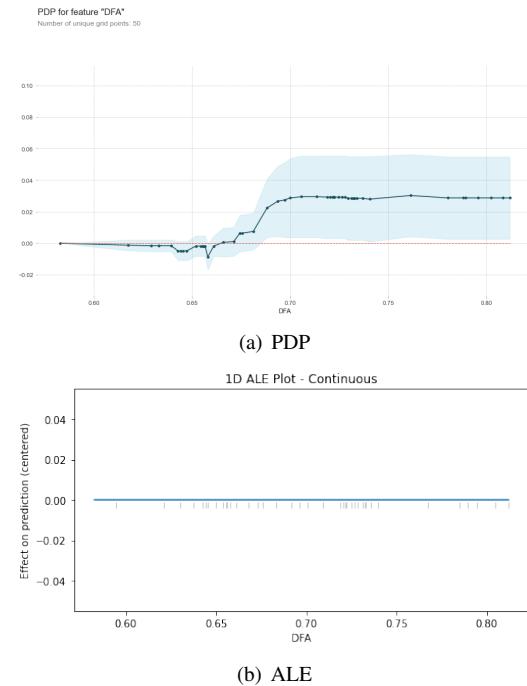


Fig. 84. PDP and ALE for the "DFA" Feature in the Random Forest

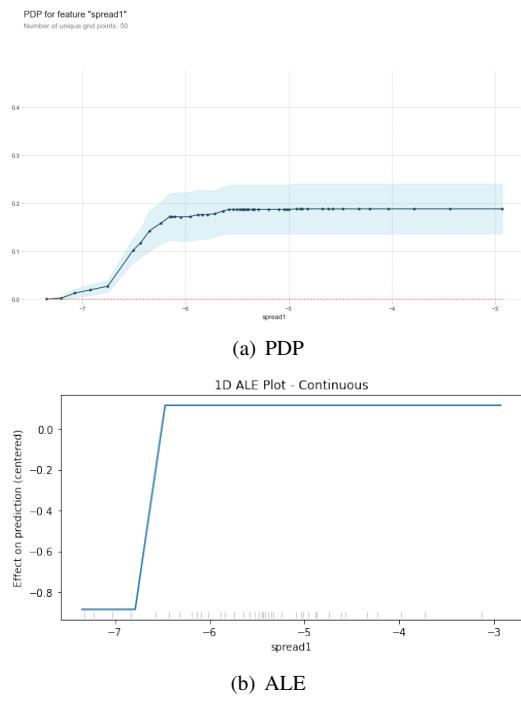


Fig. 85. PDP and ALE for the "spread1" Feature in the Random Forest

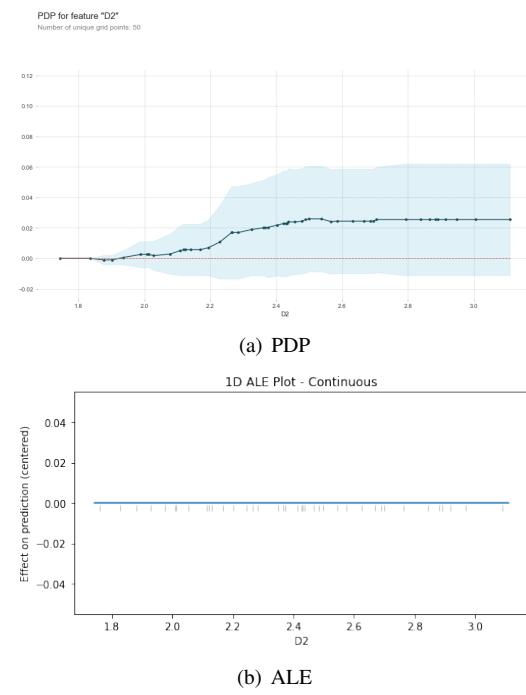


Fig. 87. PDP and ALE for the "D2" Feature in the Random Forest

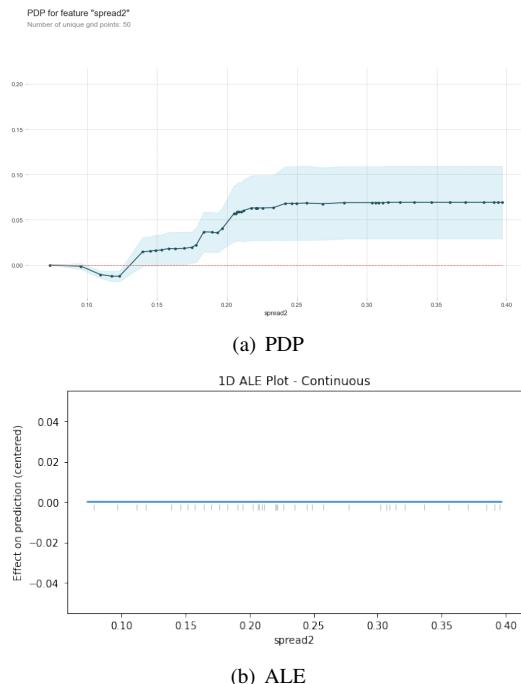


Fig. 86. PDP and ALE for the "spread2" Feature in the Random Forest

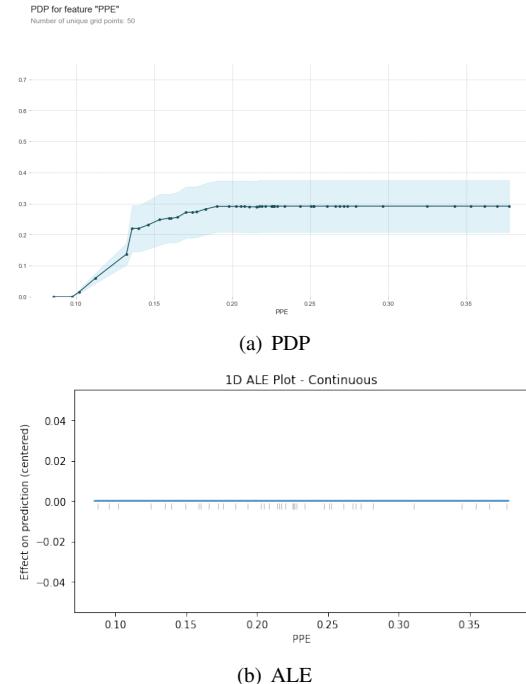
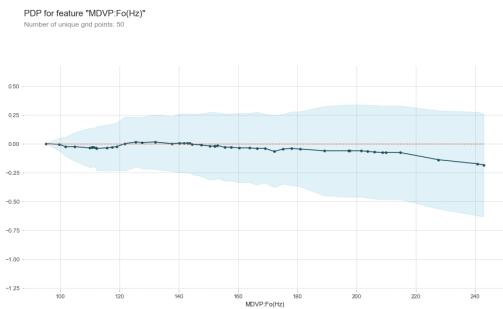
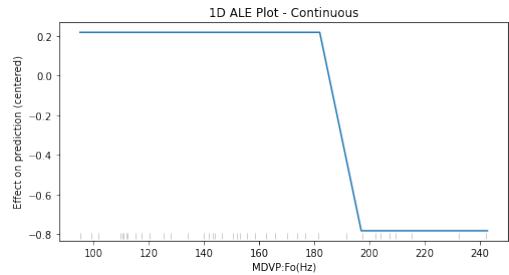


Fig. 88. PDP and ALE for the "PPE" Feature in the Random Forest

O. ALE and PDP plots for the KNN on the Parkinson

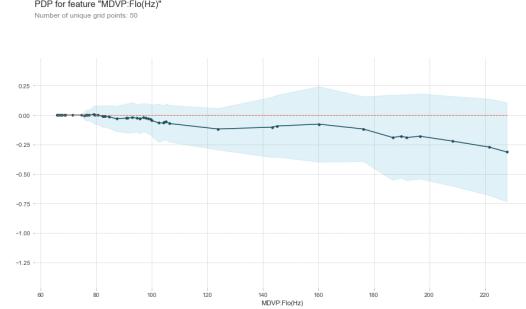


(a) PDP

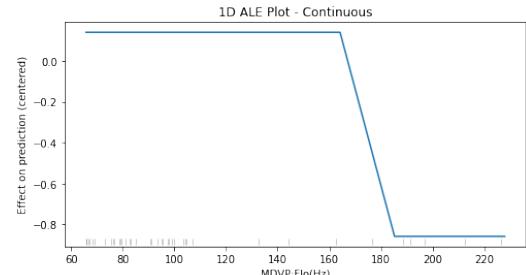


(b) ALE

Fig. 89. PDP and ALE for the "MDVP:Fo(Hz)" Feature in the knn

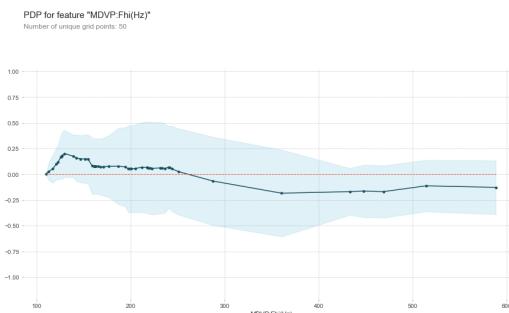


(a) PDP

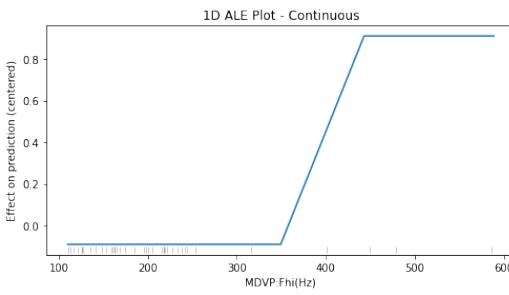


(b) ALE

Fig. 91. PDP and ALE for the "MDVP:Flo(Hz)" Feature in the knn

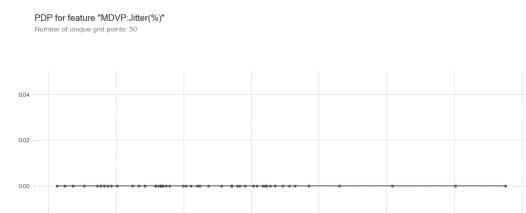


(a) PDP

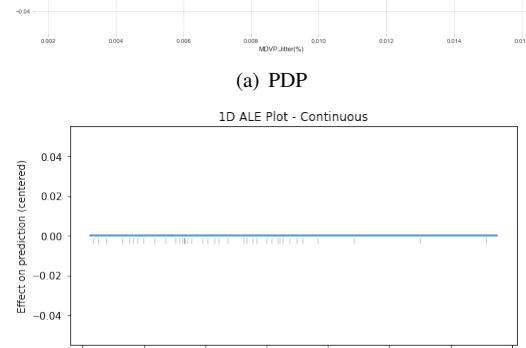


(b) ALE

Fig. 90. PDP and ALE for the "MDVP:Fhi(Hz)" Feature in the knn

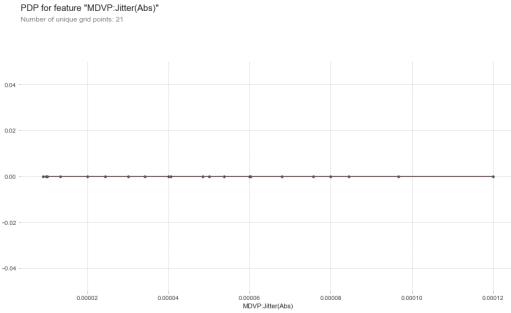


(a) PDP

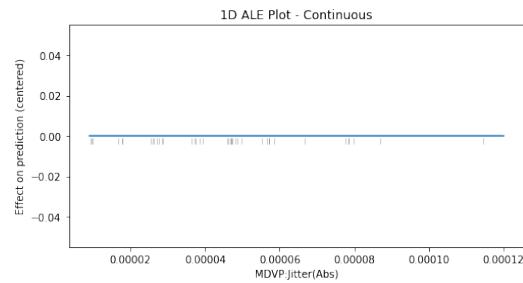


(b) ALE

Fig. 92. PDP and ALE for the "MDVP:jitter(%)" Feature in the knn

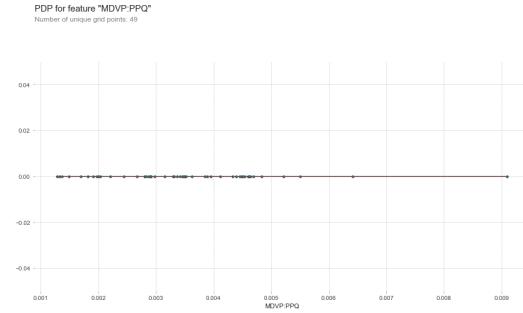


(a) PDP

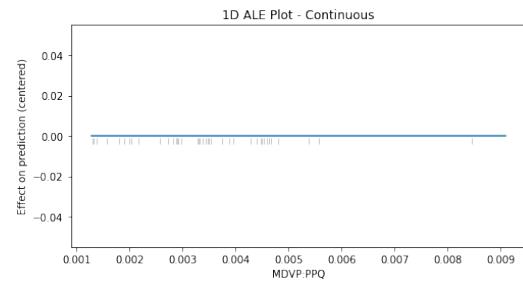


(b) ALE

Fig. 93. PDP and ALE for the "MDVP:Jitter(Abs)" Feature in the knn

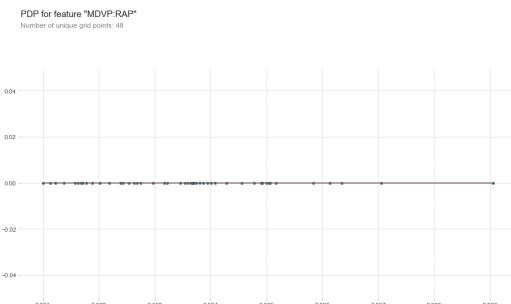


(a) PDP

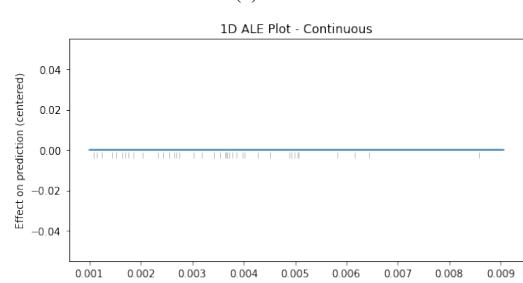


(b) ALE

Fig. 95. PDP and ALE for the "MDVP:PPQ" Feature in the knn

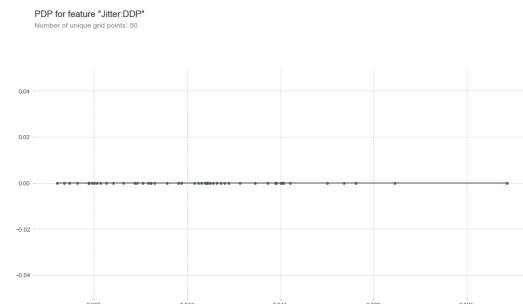


(a) PDP

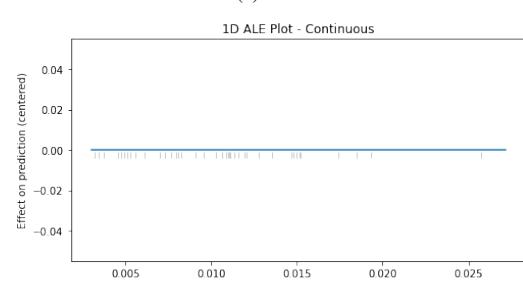


(b) ALE

Fig. 94. PDP and ALE for the "MDVP:RAP" Feature in the knn



(a) PDP



(b) ALE

Fig. 96. PDP and ALE for the "Jitter:DDP" Feature in the knn

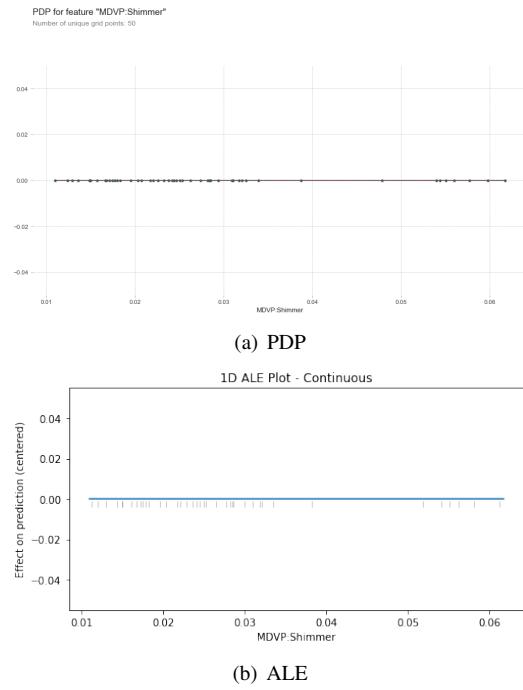


Fig. 97. PDP and ALE for the "MDVP:Shimmer" Feature in the knn

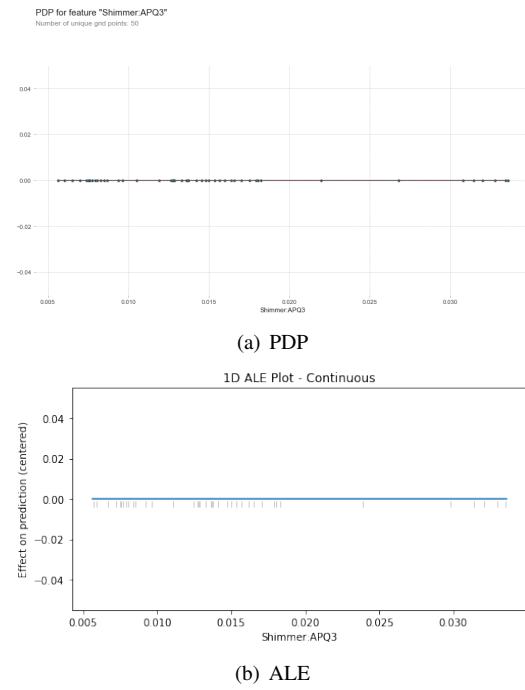


Fig. 99. PDP and ALE for the "Shimmer:APQ3" Feature in the knn

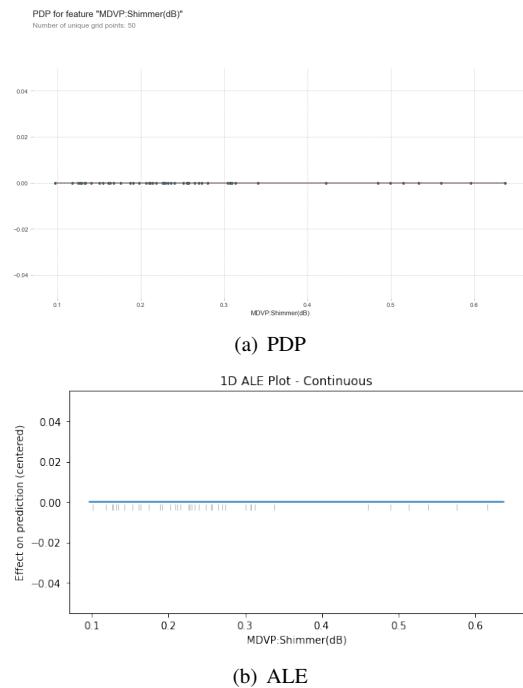


Fig. 98. PDP and ALE for the "MDVP:Shimmer(dB)" Feature in the knn

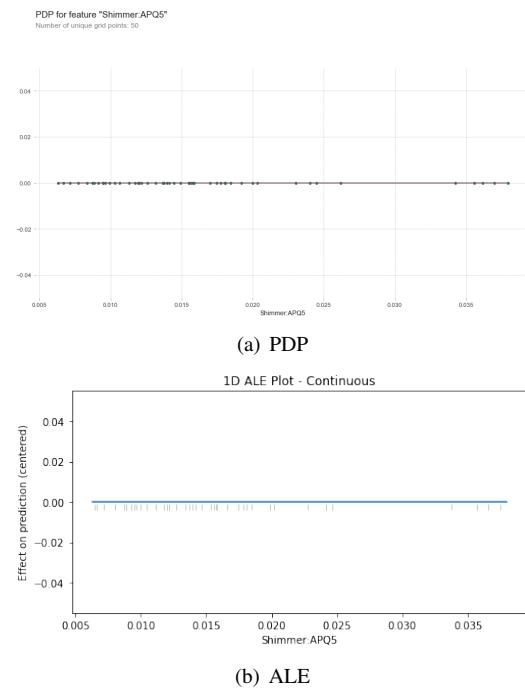
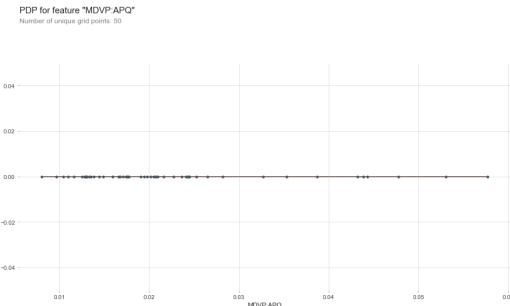
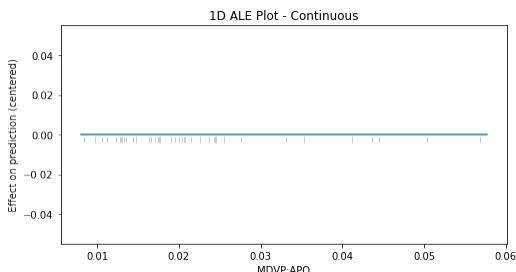


Fig. 100. PDP and ALE for the "Shimmer:APQ5" Feature in the knn

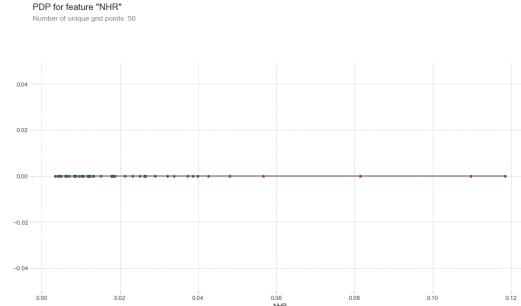


(a) PDP

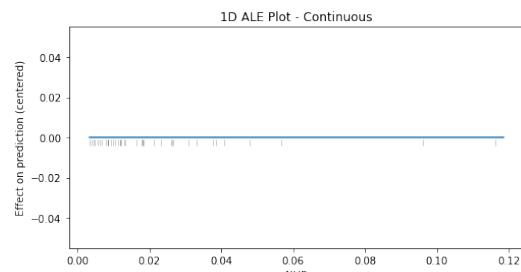


(b) ALE

Fig. 101. PDP and ALE for the "MDVP:APQ" Feature in the knn

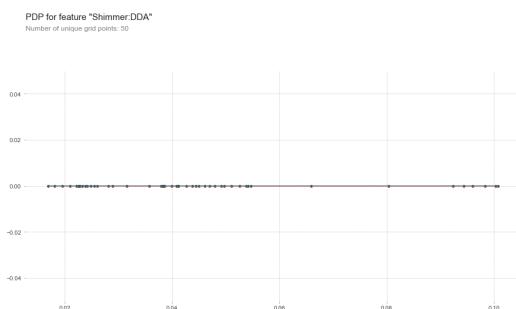


(a) PDP

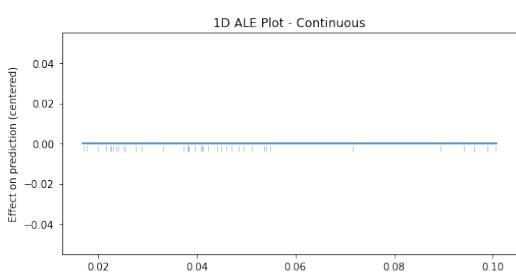


(b) ALE

Fig. 103. PDP and ALE for the "NHR" Feature in the knn



(a) PDP

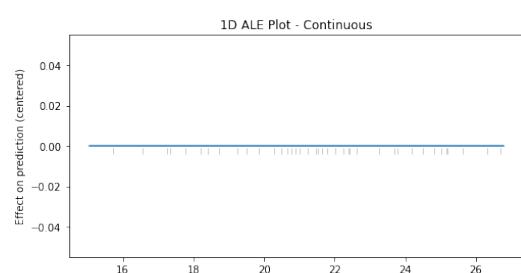


(b) ALE

Fig. 102. PDP and ALE for the "Shimmer:DDA" Feature in the knn



(a) PDP



(b) ALE

Fig. 104. PDP and ALE for the "HNR" Feature in the knn

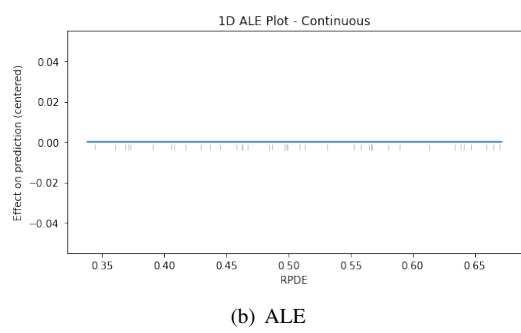
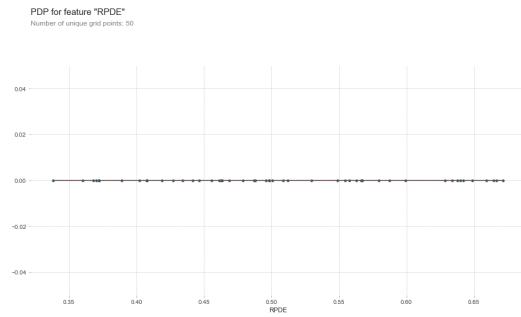


Fig. 105. PDP and ALE for the "RPDE" Feature in the knn

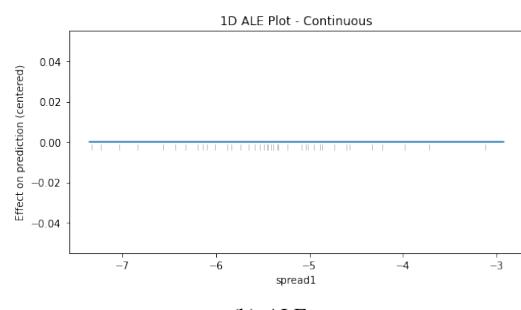
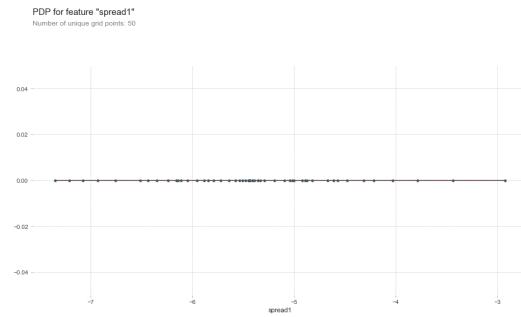


Fig. 107. PDP and ALE for the "spread1" Feature in the knn

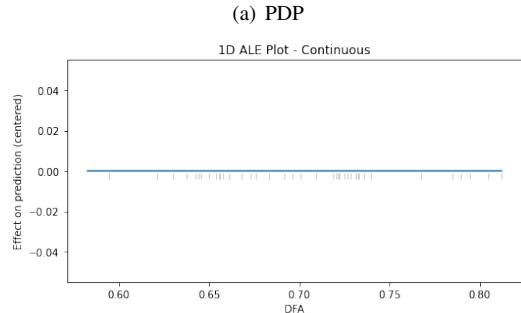
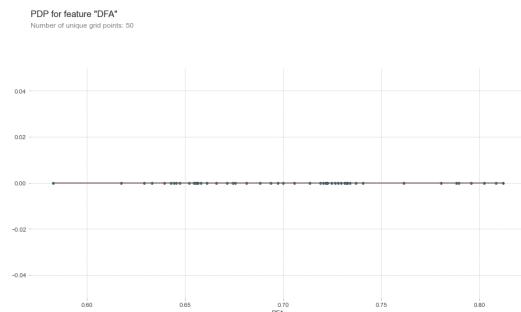


Fig. 106. PDP and ALE for the "DFA" Feature in the knn

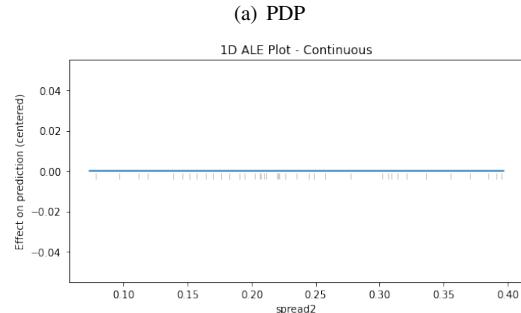
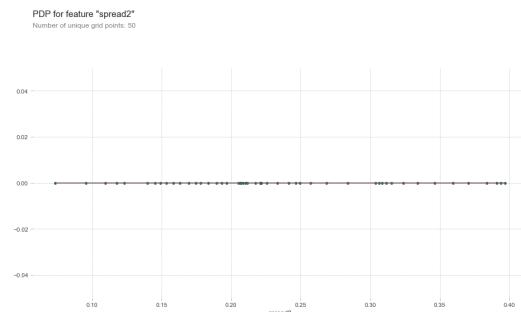


Fig. 108. PDP and ALE for the "spread2" Feature in the knn

P. ALE and PDP plots for the MLP on the Parkinson

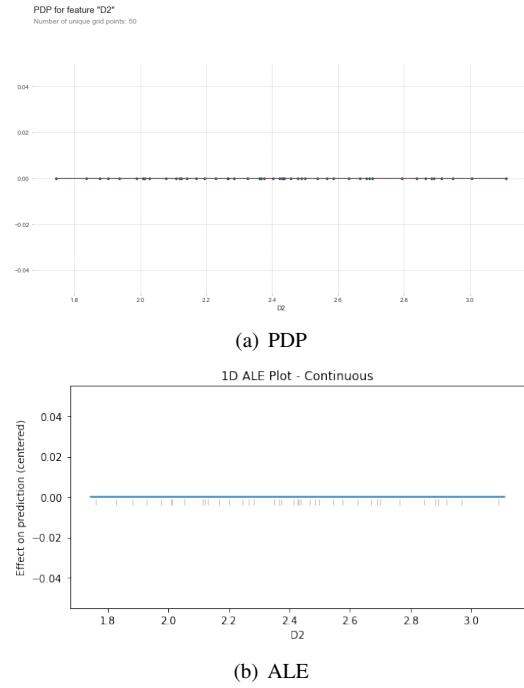


Fig. 109. PDP and ALE for the "D2" Feature in the knn

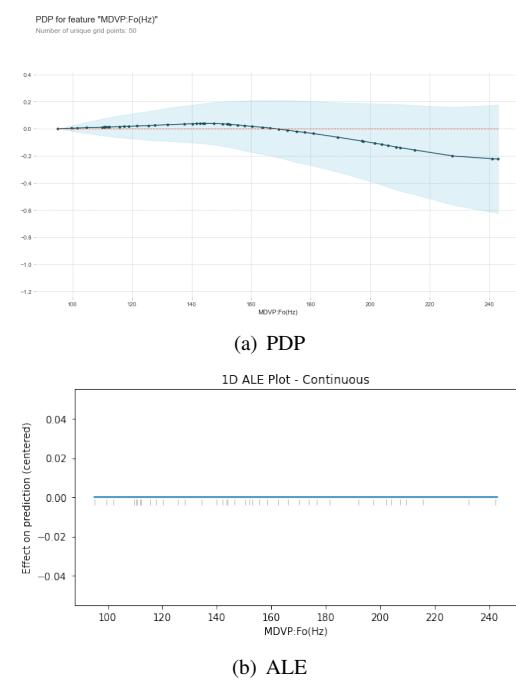


Fig. 111. PDP and ALE for the "MDVP:Fo(Hz)" Feature in the MLP

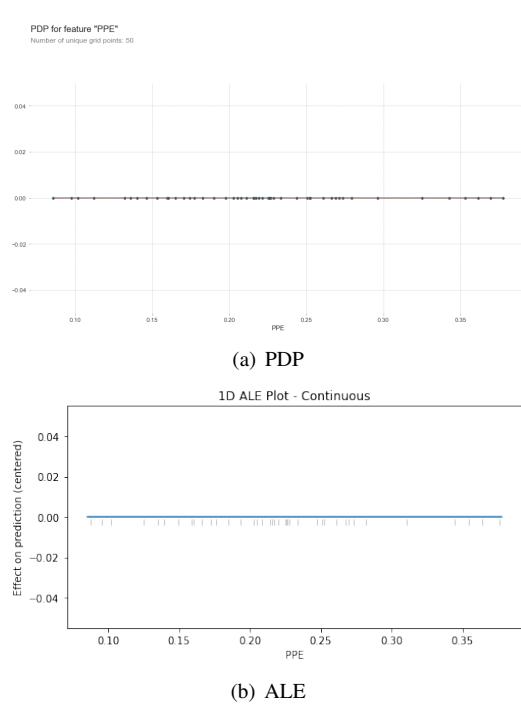


Fig. 110. PDP and ALE for the "PPE" Feature in the knn

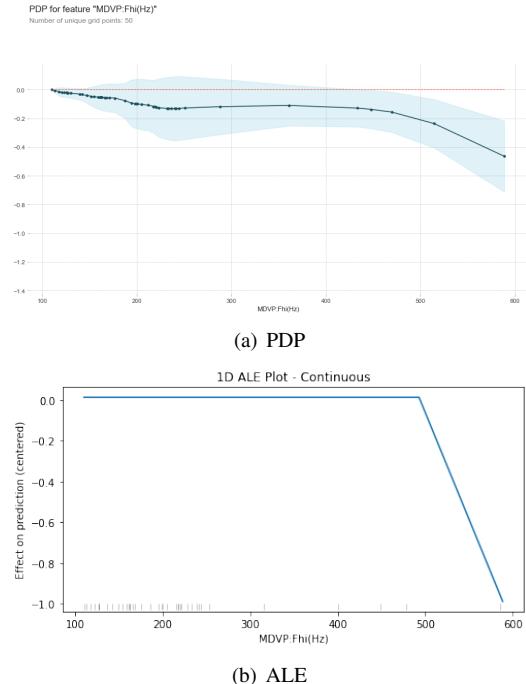
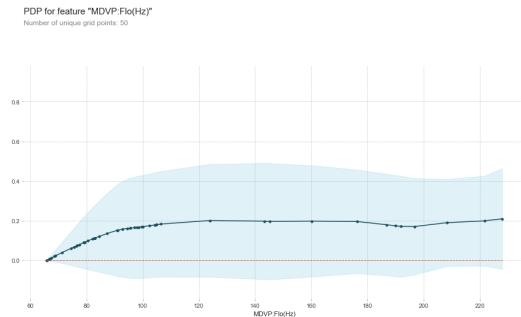
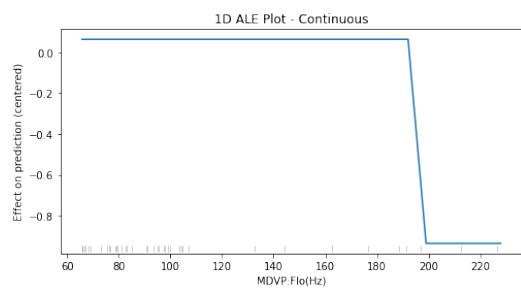


Fig. 112. PDP and ALE for the "MDVP:Fhi(Hz)" Feature in the MLP

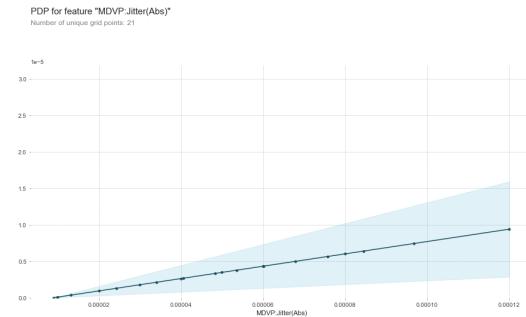


(a) PDP

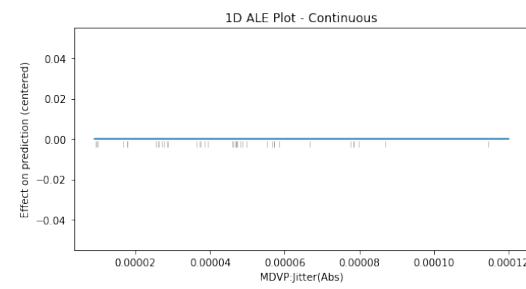


(b) ALE

Fig. 113. PDP and ALE for the "MDVP:Flo(Hz)" Feature in the MLP

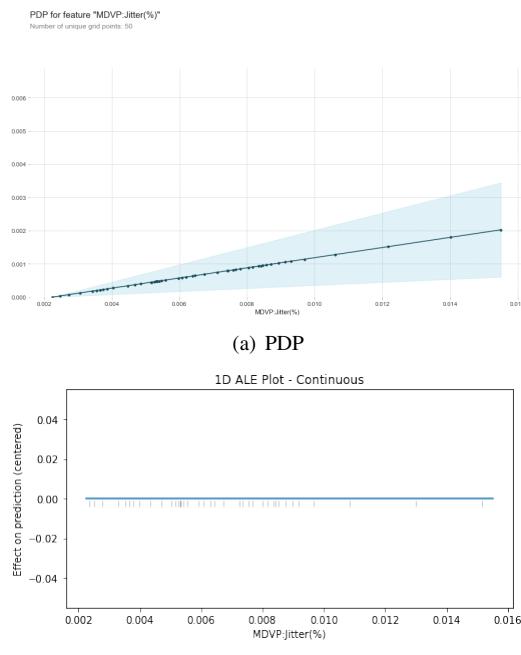


(a) PDP

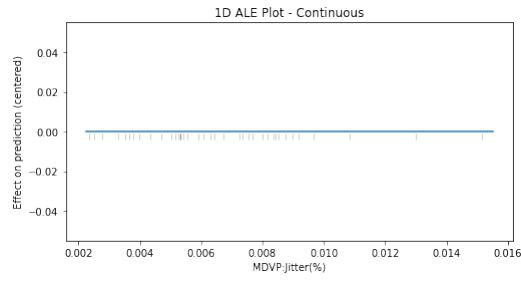


(b) ALE

Fig. 115. PDP and ALE for the "MDVP:Jitter(Abs)" Feature in the MLP

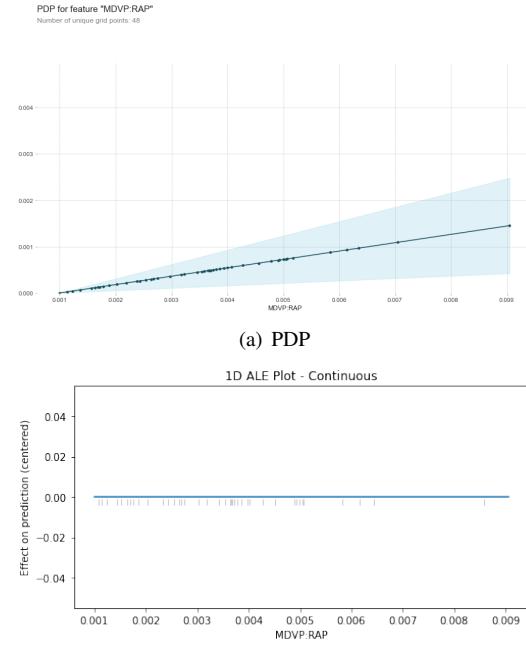


(a) PDP

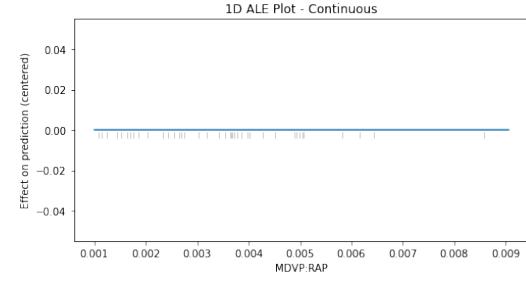


(b) ALE

Fig. 114. PDP and ALE for the "MDVP:Jitter(%)" Feature in the MLP



(a) PDP



(b) ALE

Fig. 116. PDP and ALE for the "MDVP:RAP" Feature in the MLP

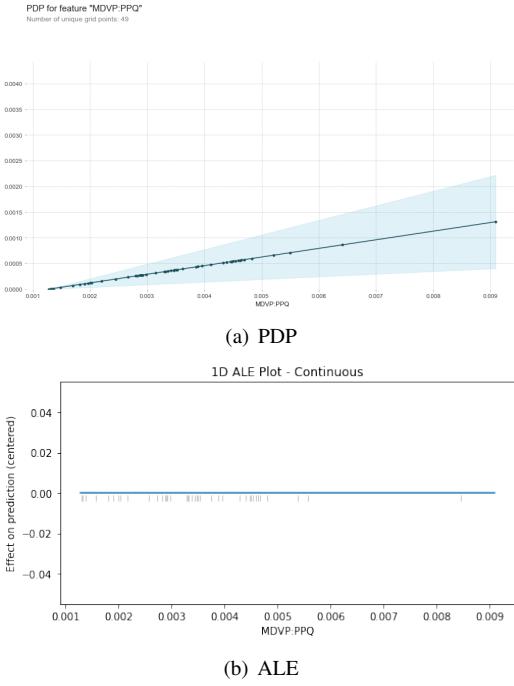


Fig. 117. PDP and ALE for the "MDVP:PPQ" Feature in the MLP

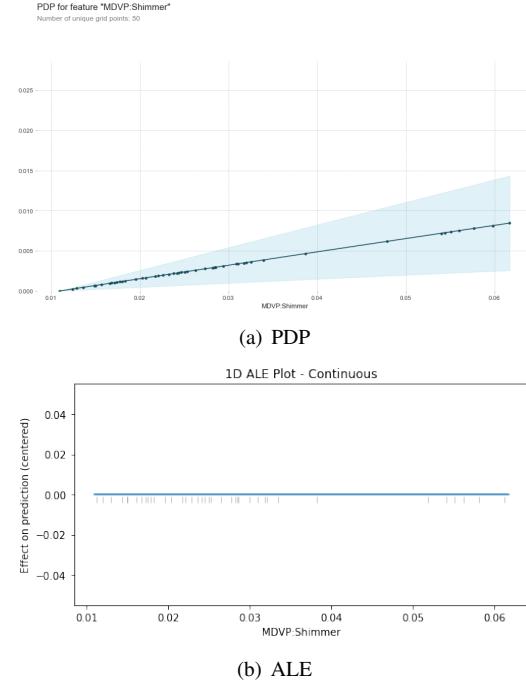


Fig. 119. PDP and ALE for the "MDVP:Shimmer" Feature in the MLP

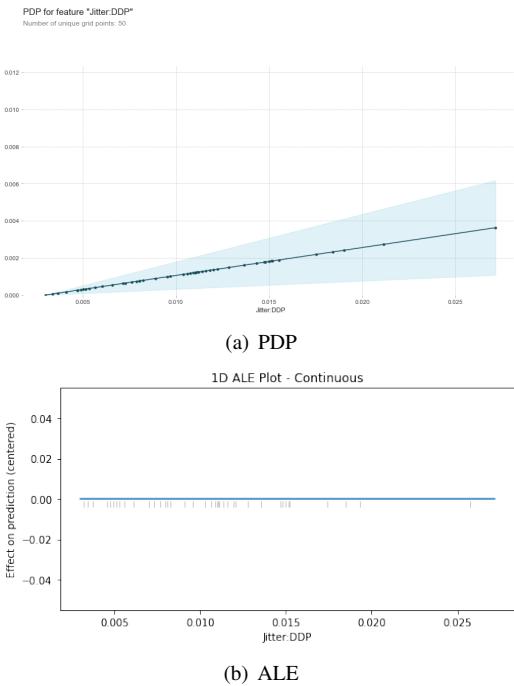


Fig. 118. PDP and ALE for the "Jitter:DDP" Feature in the MLP

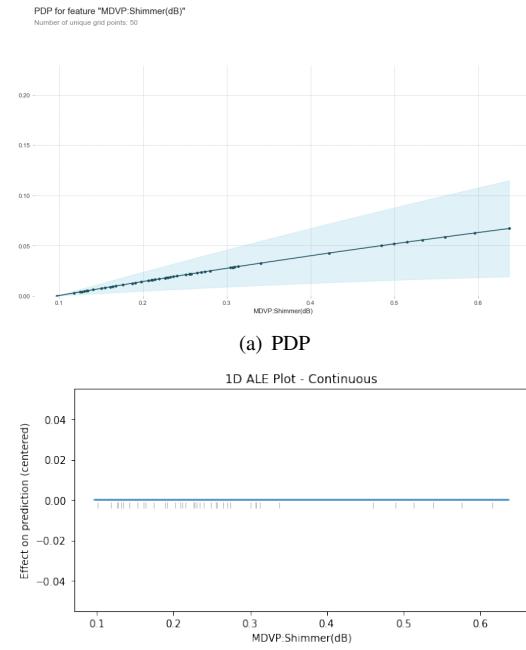


Fig. 120. PDP and ALE for the "MDVP:Shimmer(dB)" Feature in the MLP

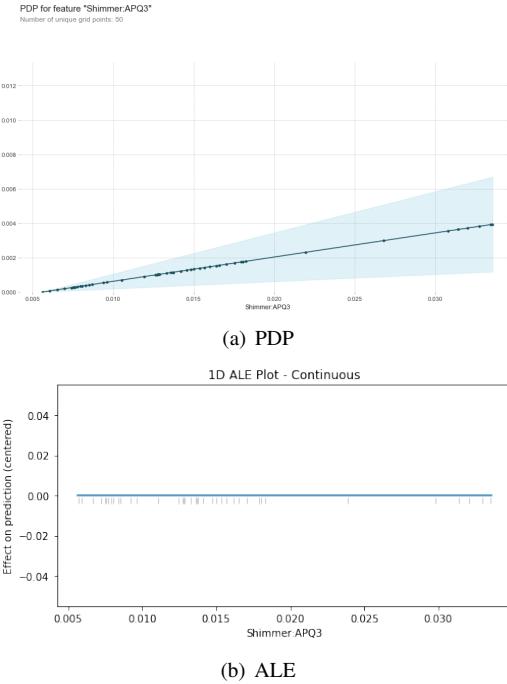


Fig. 121. PDP and ALE for the "Shimmer:APQ3" Feature in the MLP

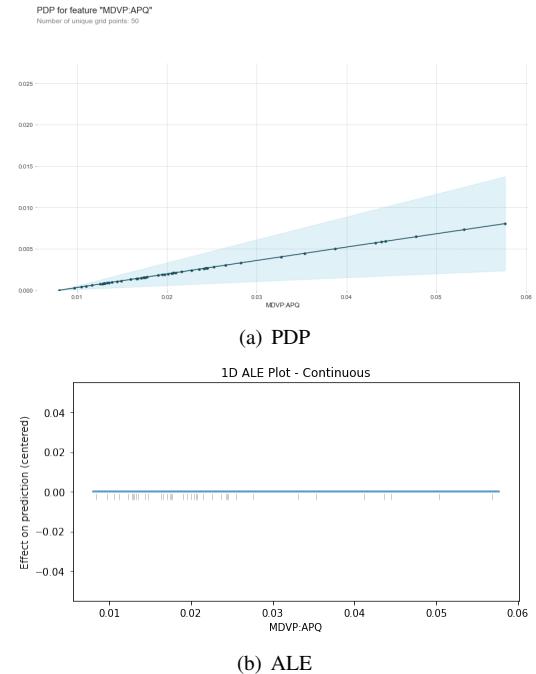


Fig. 123. PDP and ALE for the "MDVP:APQ" Feature in the MLP

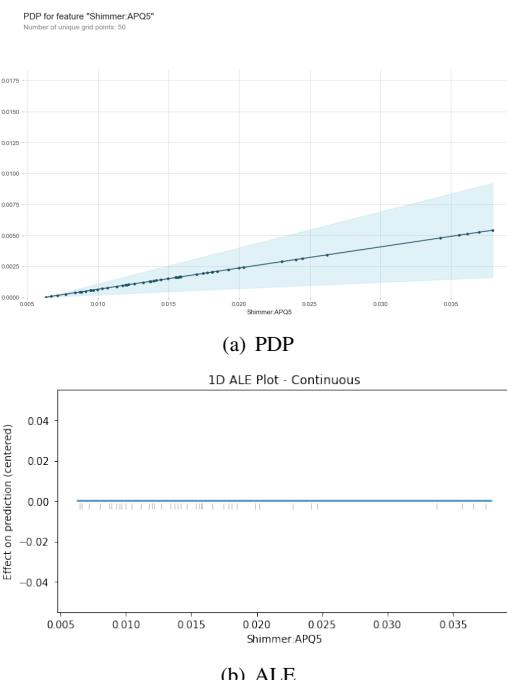


Fig. 122. PDP and ALE for the "Shimmer:APQ5" Feature in the MLP

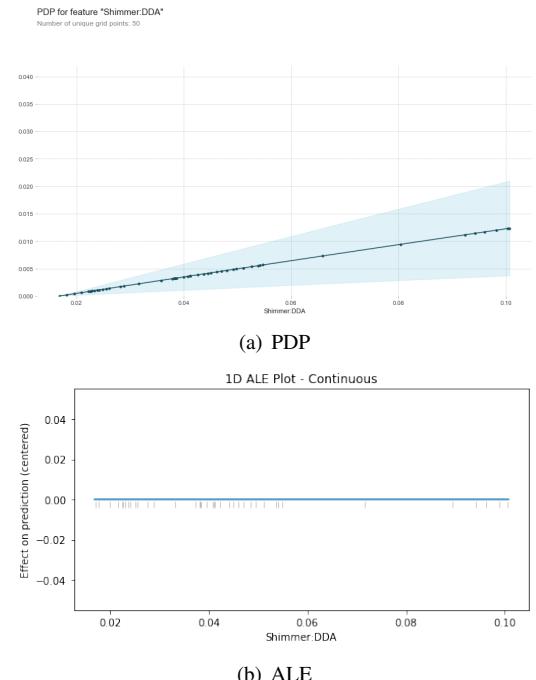


Fig. 124. PDP and ALE for the "Shimmer:DDA" Feature in the MLP

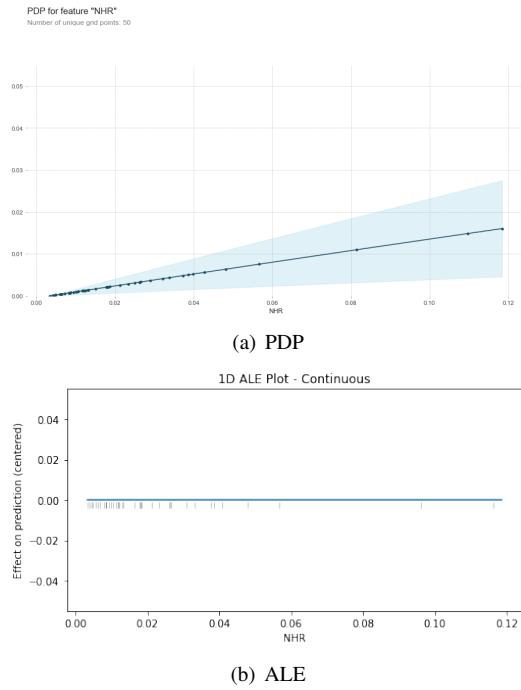


Fig. 125. PDP and ALE for the "NHR" Feature in the MLP

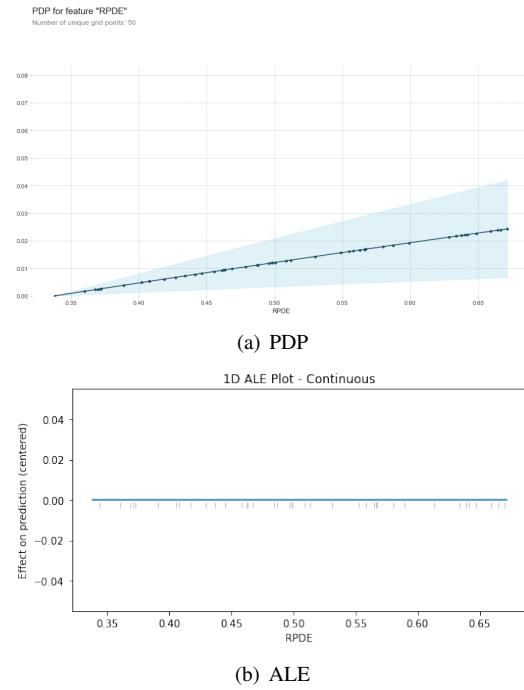


Fig. 127. PDP and ALE for the "RPDE" Feature in the MLP

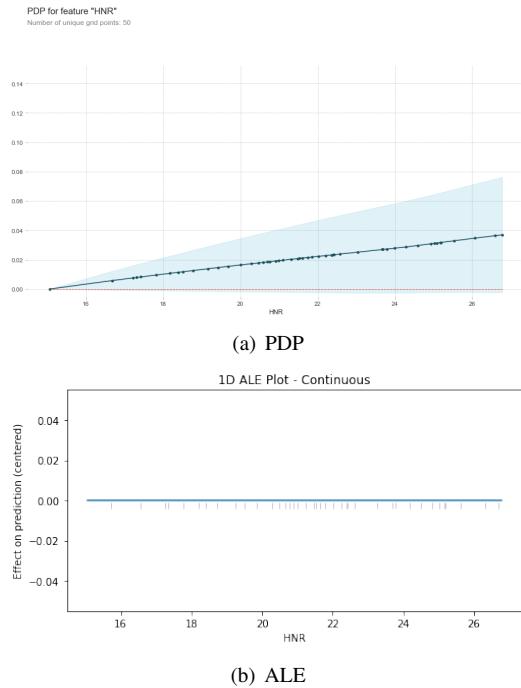


Fig. 126. PDP and ALE for the "HNR" Feature in the MLP

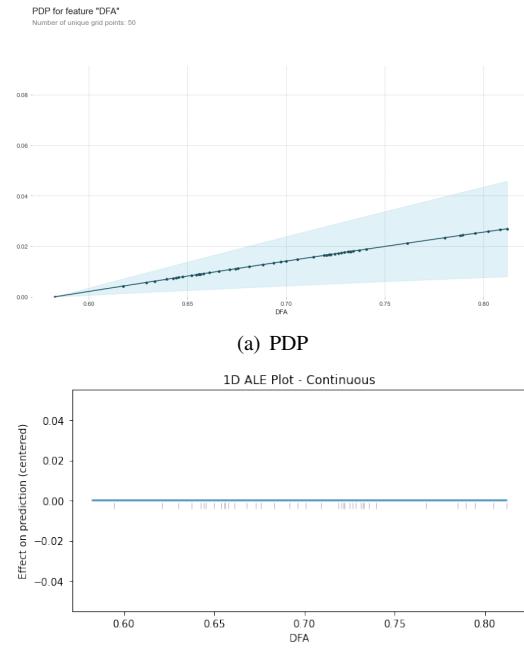


Fig. 128. PDP and ALE for the "DFA" Feature in the MLP

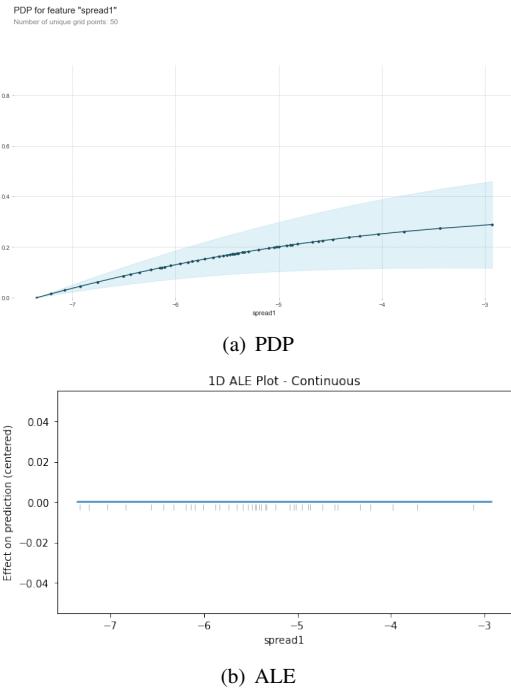


Fig. 129. PDP and ALE for the "spread1" Feature in the MLP

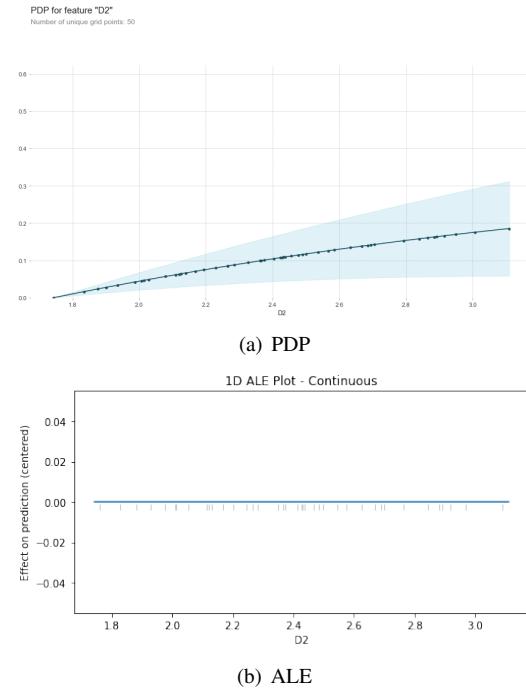


Fig. 131. PDP and ALE for the "D2" Feature in the MLP

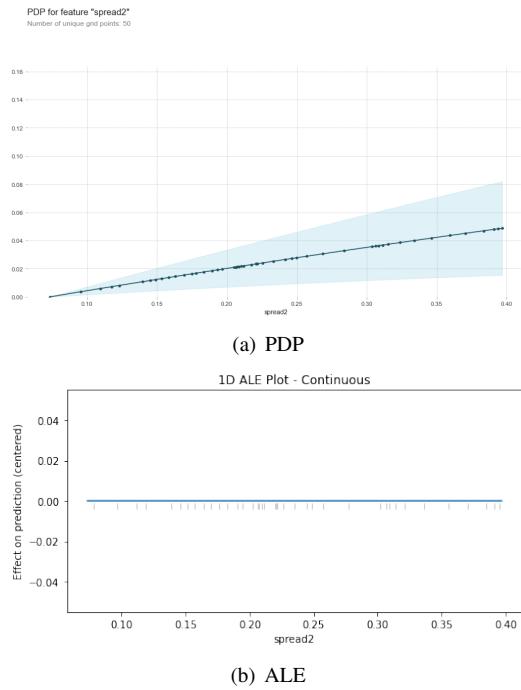


Fig. 130. PDP and ALE for the "spread2" Feature in the MLP

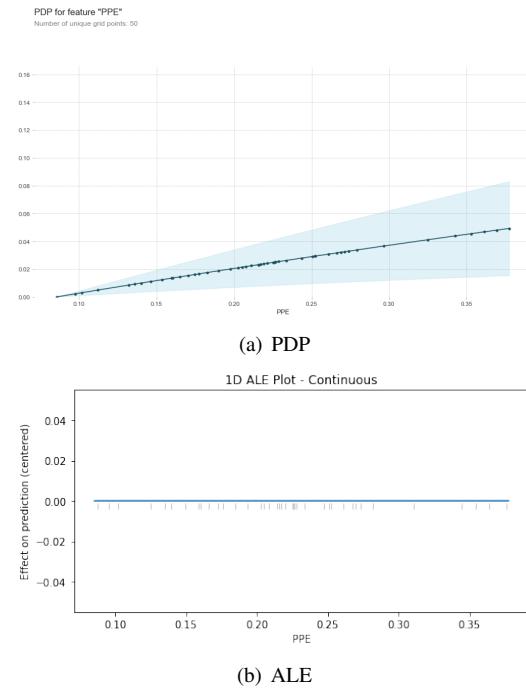
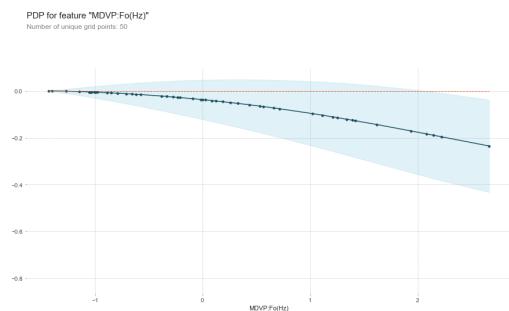
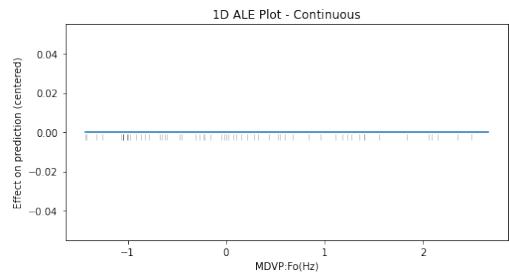


Fig. 132. PDP and ALE for the "PPE" Feature in the MLP

Q. ALE and PDP plots for the SVM on the Parkinson

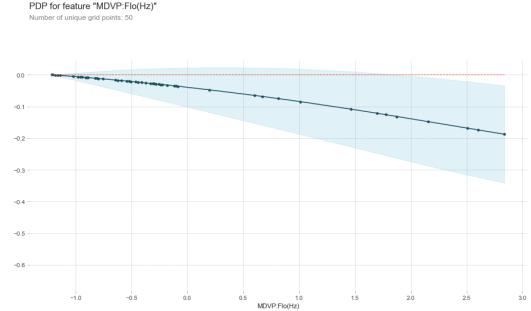


(a) PDP

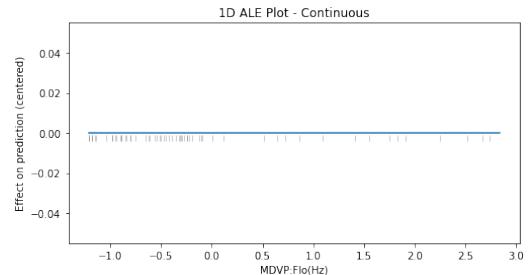


(b) ALE

Fig. 133. PDP and ALE for the "MDVP:Fo(Hz)" Feature in the SVM

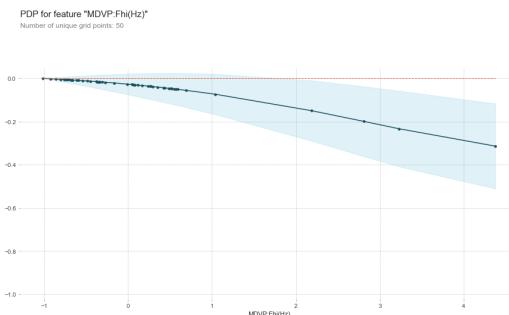


(a) PDP

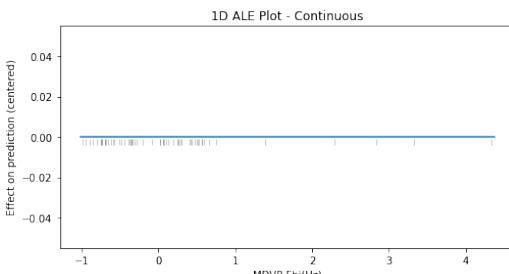


(b) ALE

Fig. 135. PDP and ALE for the "MDVP:Flo(Hz)" Feature in the SVM

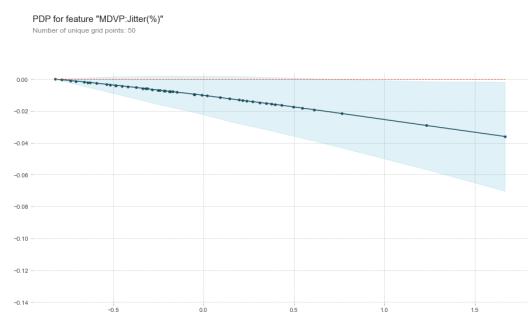


(a) PDP

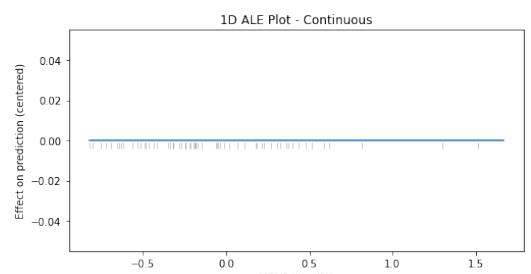


(b) ALE

Fig. 134. PDP and ALE for the "MDVP:Fhi(Hz)" Feature in the SVM

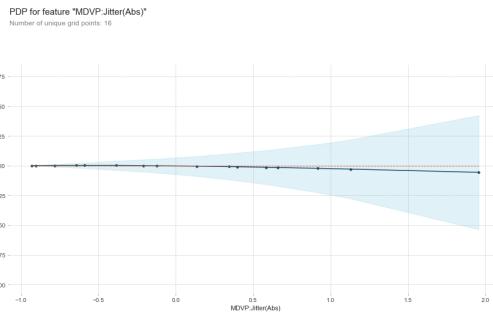


(a) PDP

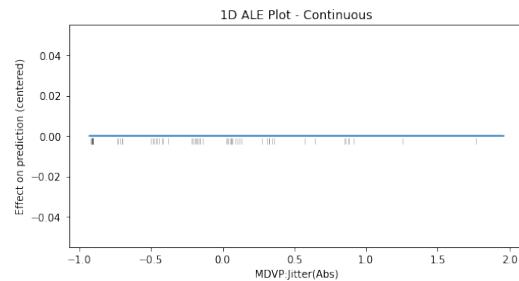


(b) ALE

Fig. 136. PDP and ALE for the "MDVP:Jitter(%)" Feature in the SVM

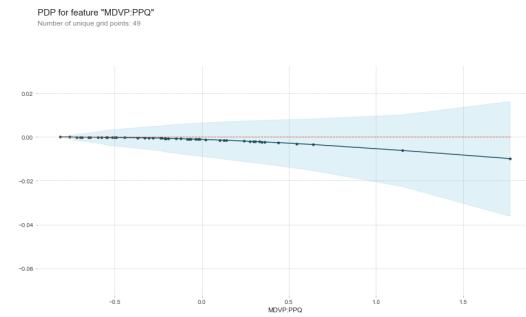


(a) PDP

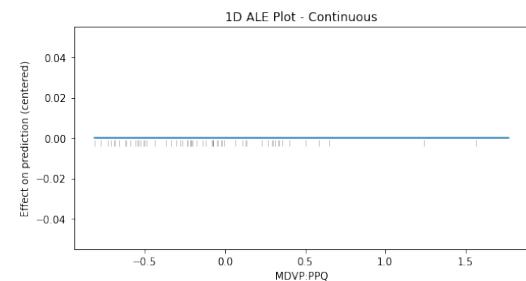


(b) ALE

Fig. 137. PDP and ALE for the "MDVP:Jitter(Abs)" Feature in the SVM

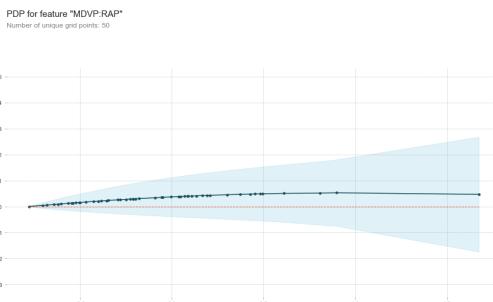


(a) PDP

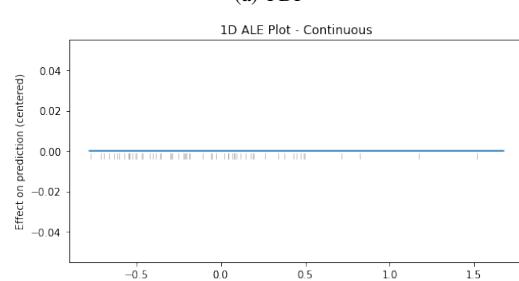


(b) ALE

Fig. 139. PDP and ALE for the "MDVP:PPQ" Feature in the SVM

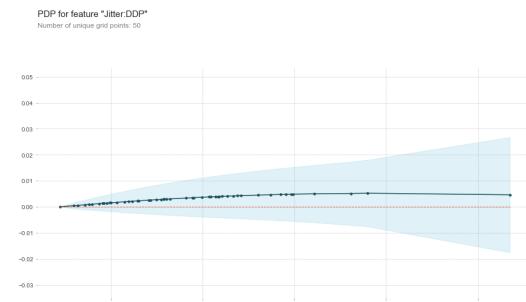


(a) PDP

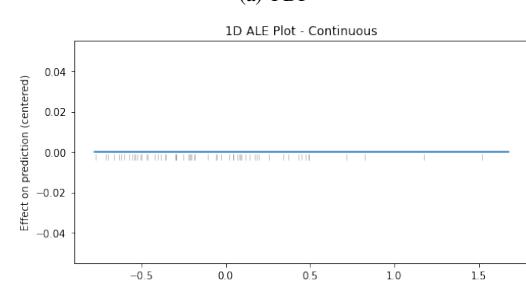


(b) ALE

Fig. 138. PDP and ALE for the "MDVP:RAP" Feature in the SVM

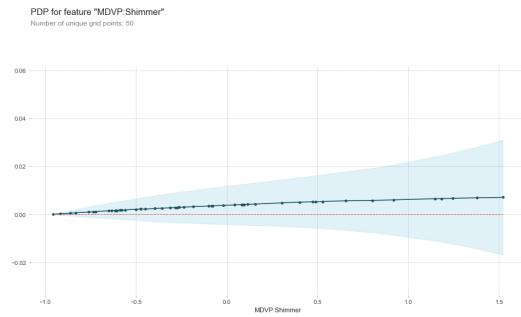


(a) PDP

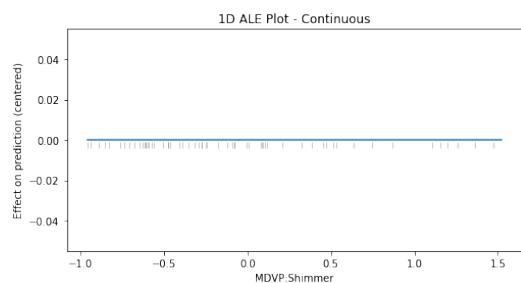


(b) ALE

Fig. 140. PDP and ALE for the "Jitter:DDP" Feature in the SVM

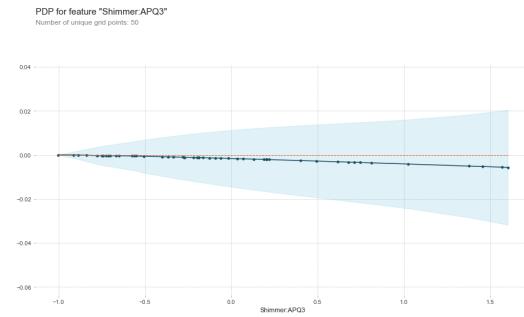


(a) PDP

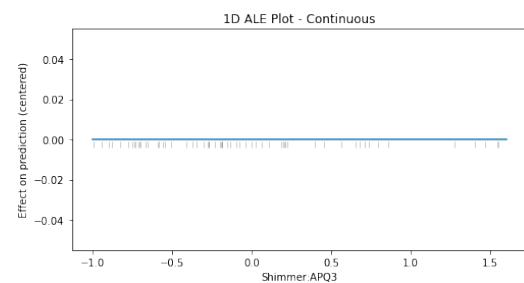


(b) ALE

Fig. 141. PDP and ALE for the "MDVP:Shimmer" Feature in the SVM

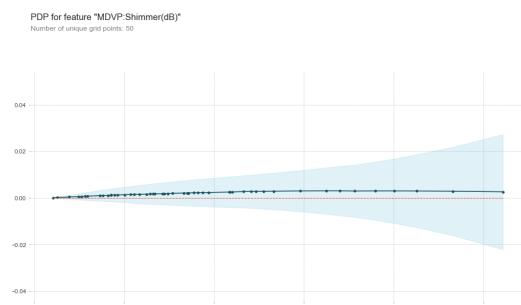


(a) PDP

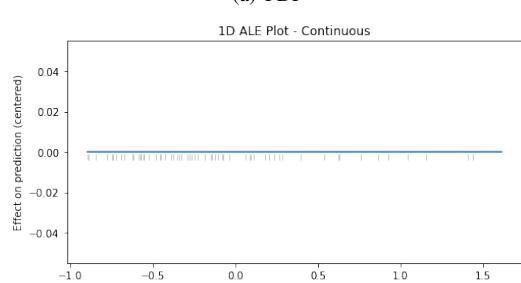


(b) ALE

Fig. 143. PDP and ALE for the "Shimmer:APQ3" Feature in the SVM

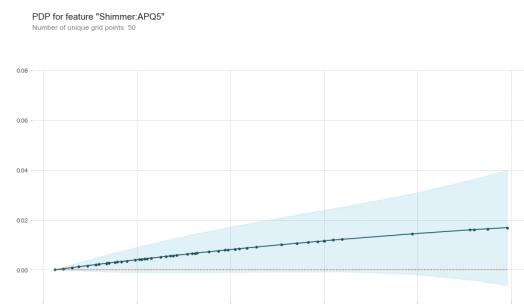


(a) PDP

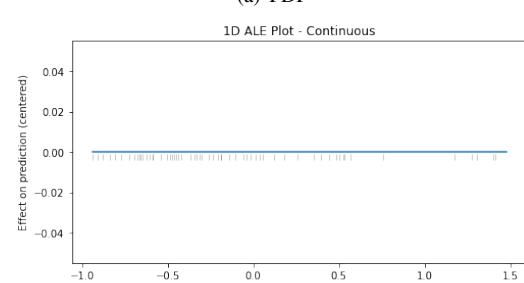


(b) ALE

Fig. 142. PDP and ALE for the "MDVP:Shimmer(dB)" Feature in the SVM

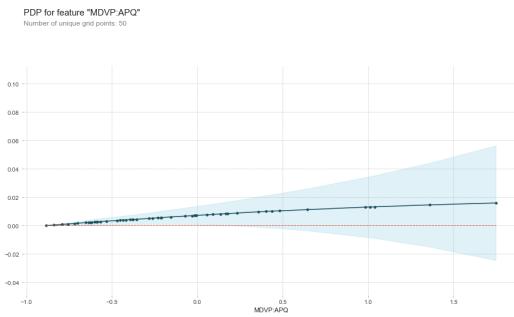


(a) PDP

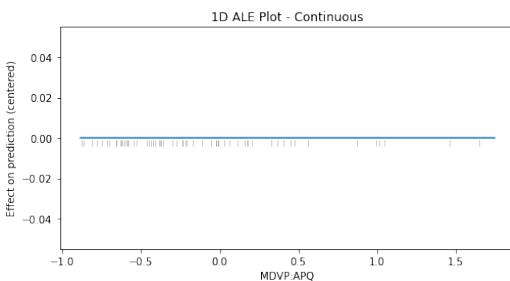


(b) ALE

Fig. 144. PDP and ALE for the "Shimmer:APQ5" Feature in the SVM

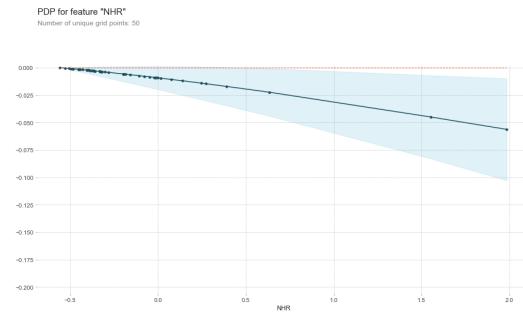


(a) PDP

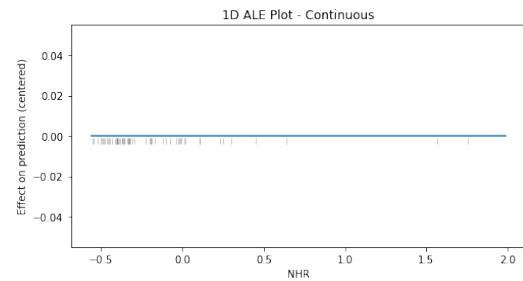


(b) ALE

Fig. 145. PDP and ALE for the "MDVP:APQ" Feature in the SVM

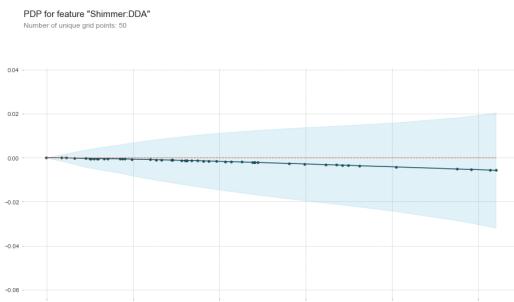


(a) PDP

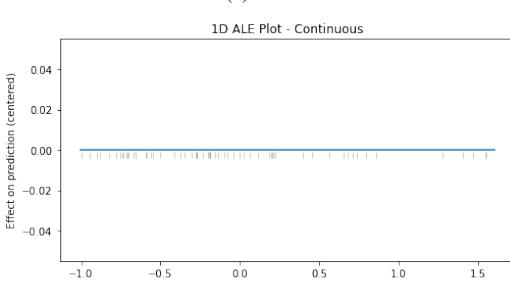


(b) ALE

Fig. 147. PDP and ALE for the "NHR" Feature in the SVM

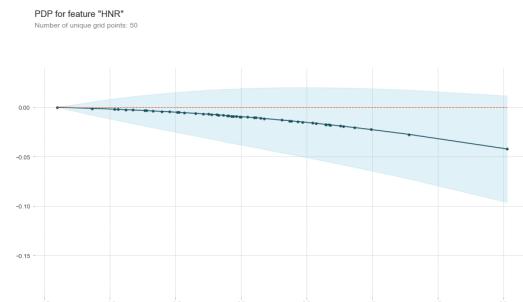


(a) PDP

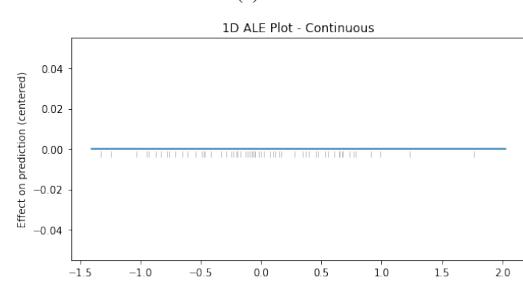


(b) ALE

Fig. 146. PDP and ALE for the "Shimmer:DDA" Feature in the SVM

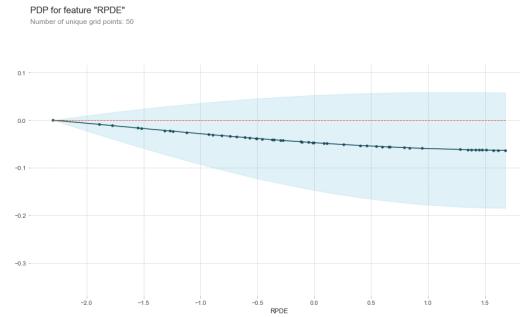


(a) PDP

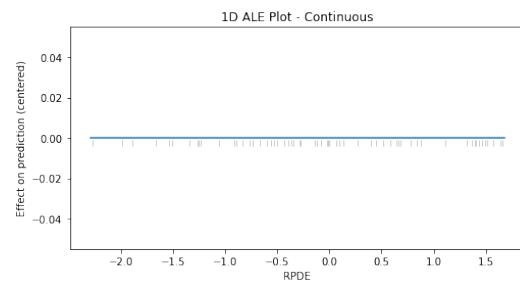


(b) ALE

Fig. 148. PDP and ALE for the "HNR" Feature in the SVM

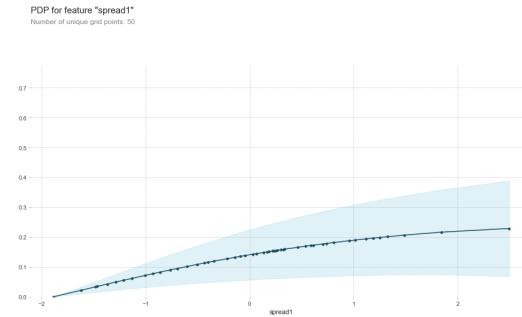


(a) PDP

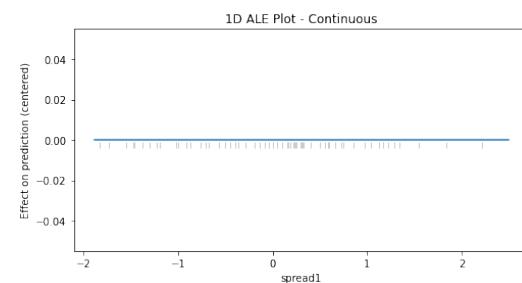


(b) ALE

Fig. 149. PDP and ALE for the "RPDE" Feature in the SVM

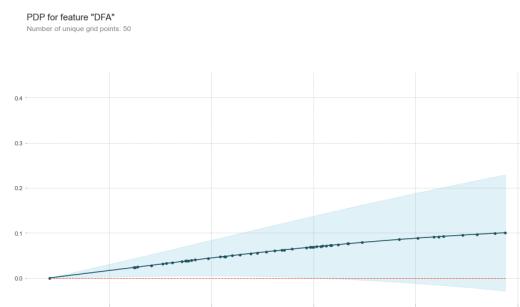


(a) PDP

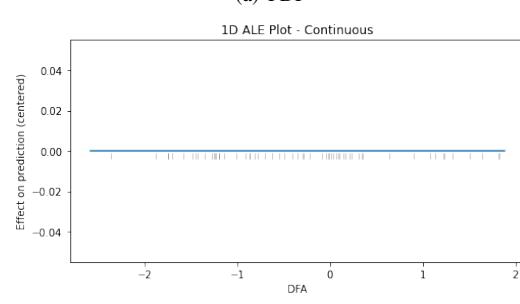


(b) ALE

Fig. 151. PDP and ALE for the "spread1" Feature in the SVM

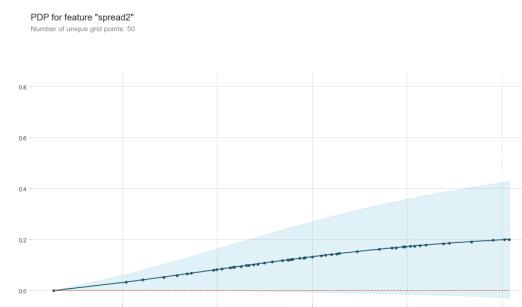


(a) PDP

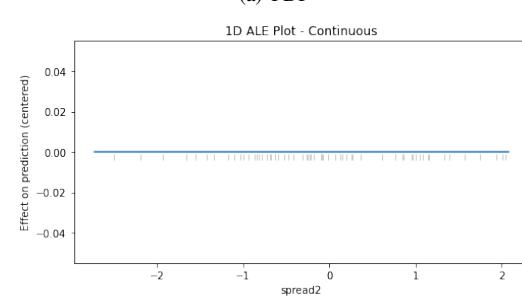


(b) ALE

Fig. 150. PDP and ALE for the "DFA" Feature in the SVM



(a) PDP



(b) ALE

Fig. 152. PDP and ALE for the "spread2" Feature in the SVM

R. ALE and PDP plots for the K-Means based classifier on the Parkinson

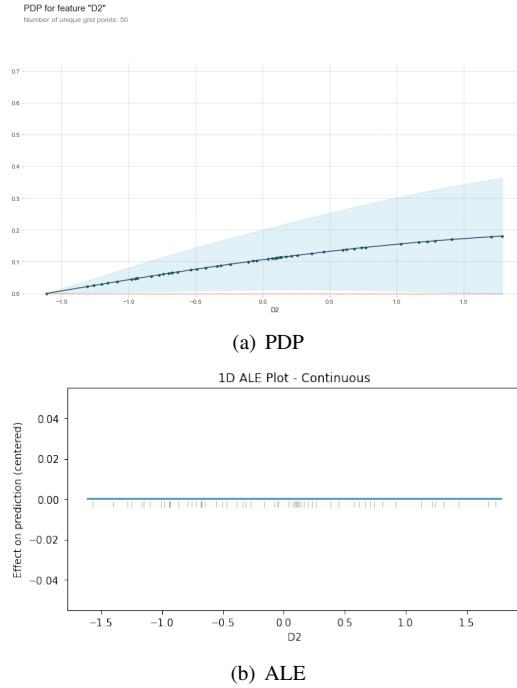


Fig. 153. PDP and ALE for the "D2" Feature in the SVM

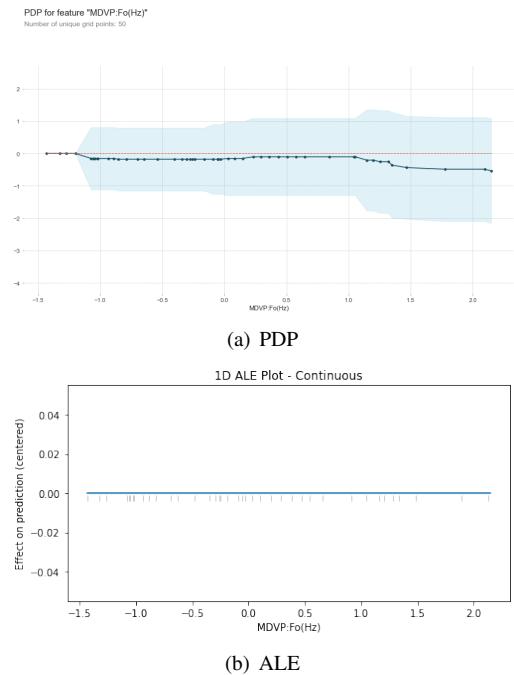


Fig. 155. PDP and ALE for the "MDVP:Fo(Hz)" Feature in the K-Means based classifier

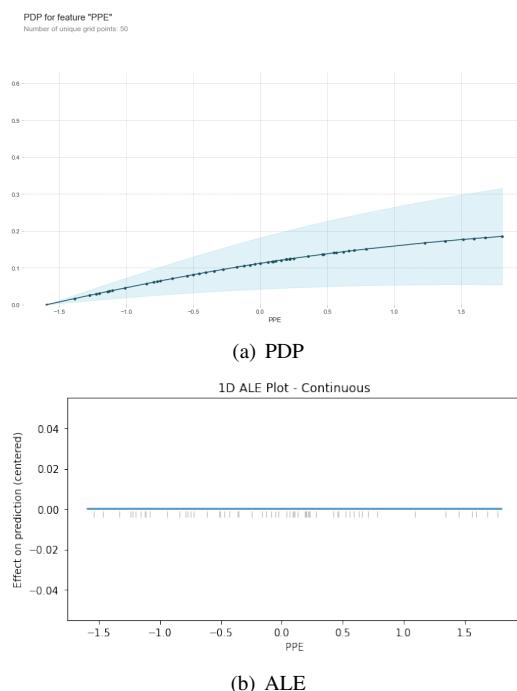
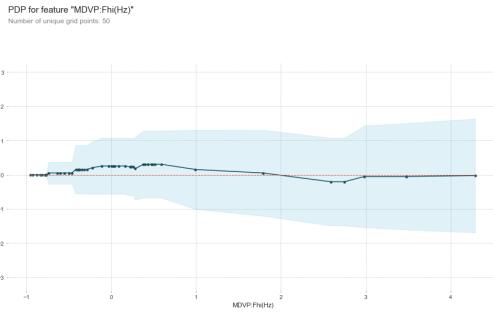
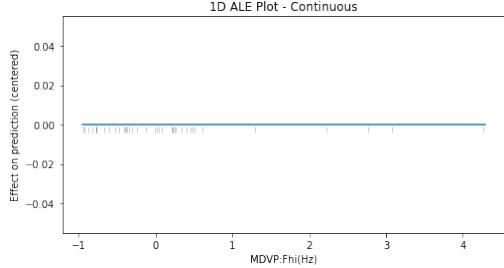


Fig. 154. PDP and ALE for the "PPE" Feature in the SVM

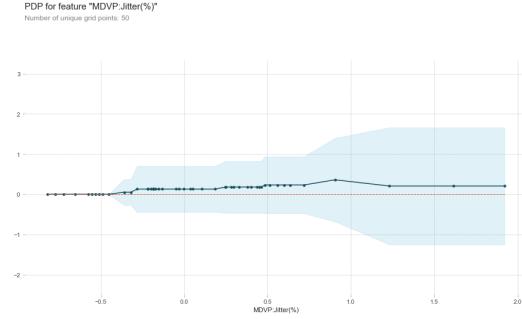


(a) PDP

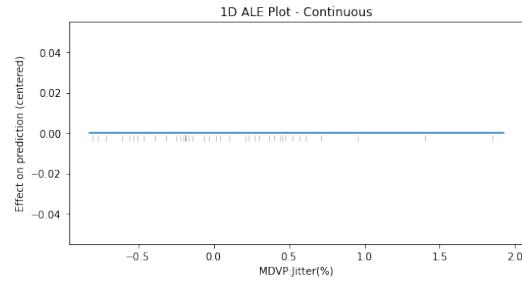


(b) ALE

Fig. 156. PDP and ALE for the "MDVP:Fhi(Hz)" Feature in the K-Means based classifier

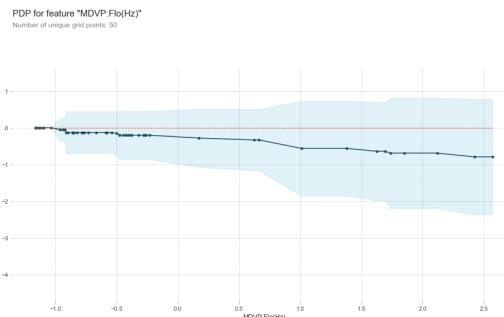


(a) PDP

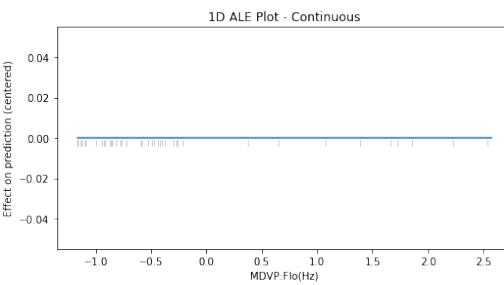


(b) ALE

Fig. 158. PDP and ALE for the "MDVP:Jitter(%)" Feature in the K-Means based classifier

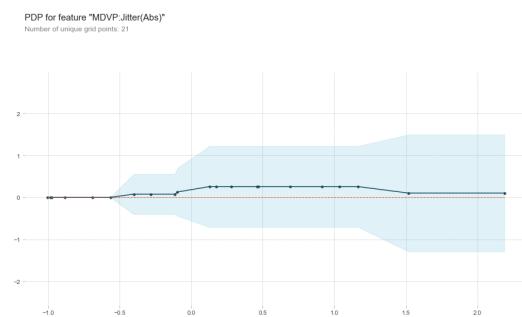


(a) PDP

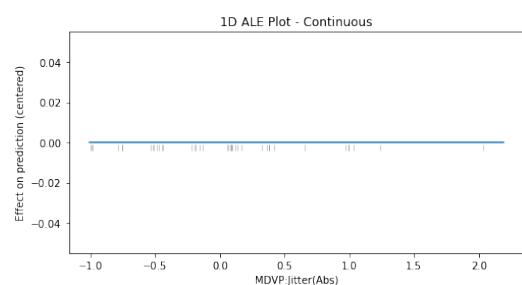


(b) ALE

Fig. 157. PDP and ALE for the "MDVP:Flo(Hz)" Feature in the K-Means based classifier

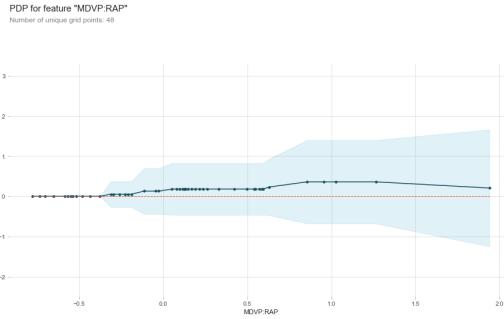


(a) PDP

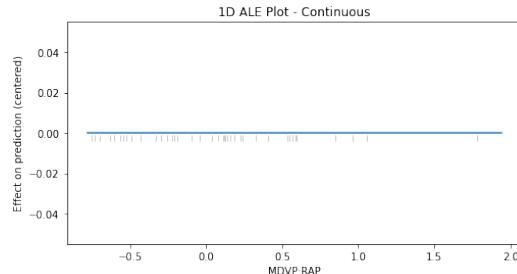


(b) ALE

Fig. 159. PDP and ALE for the "MDVP:Jitter(Abs)" Feature in the K-Means based classifier

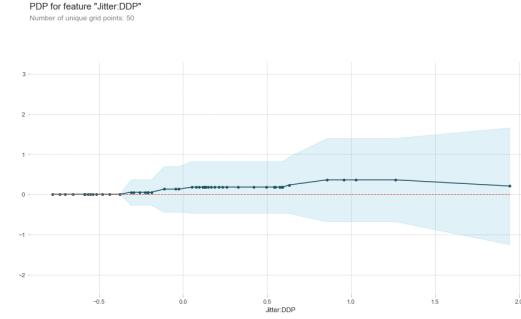


(a) PDP

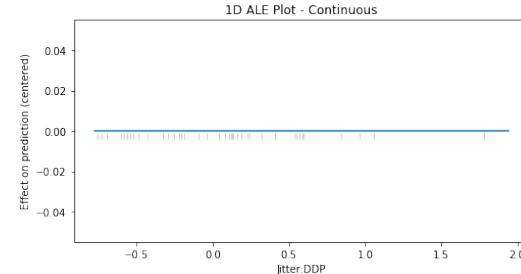


(b) ALE

Fig. 160. PDP and ALE for the "MDVP:RAP" Feature in the K-Means based classifier

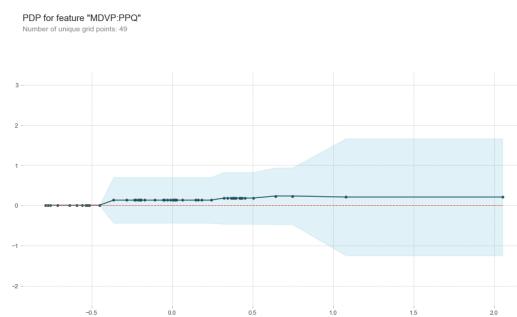


(a) PDP

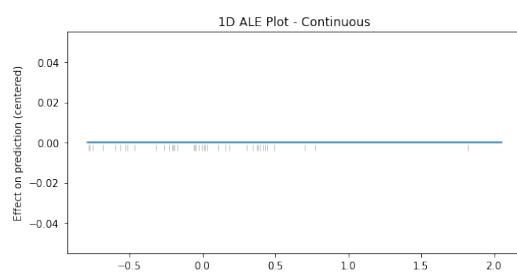


(b) ALE

Fig. 162. PDP and ALE for the "Jitter:DDP" Feature in the K-Means based classifier

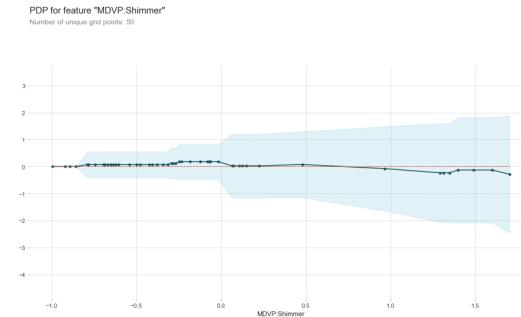


(a) PDP

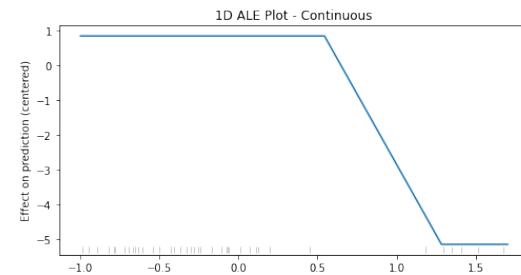


(b) ALE

Fig. 161. PDP and ALE for the "MDVP:PPQ" Feature in the K-Means based classifier



(a) PDP



(b) ALE

Fig. 163. PDP and ALE for the "MDVP:Shimmer" Feature in the K-Means based classifier

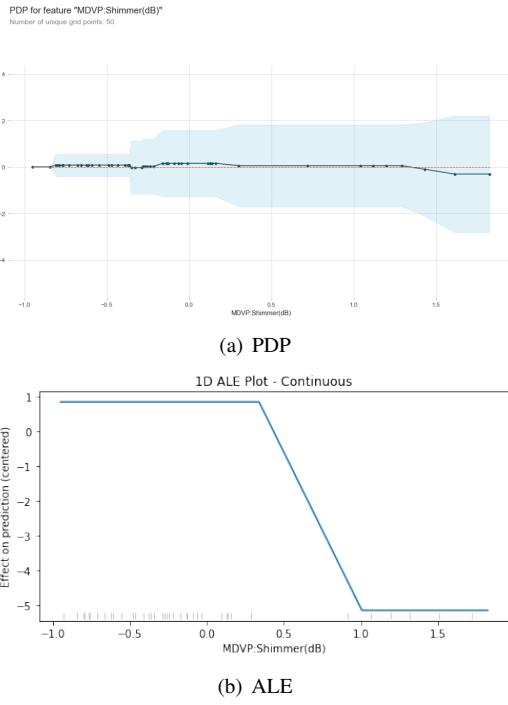


Fig. 164. PDP and ALE for the "MDVP:Shimmer(dB)" Feature in the K-Means based classifier

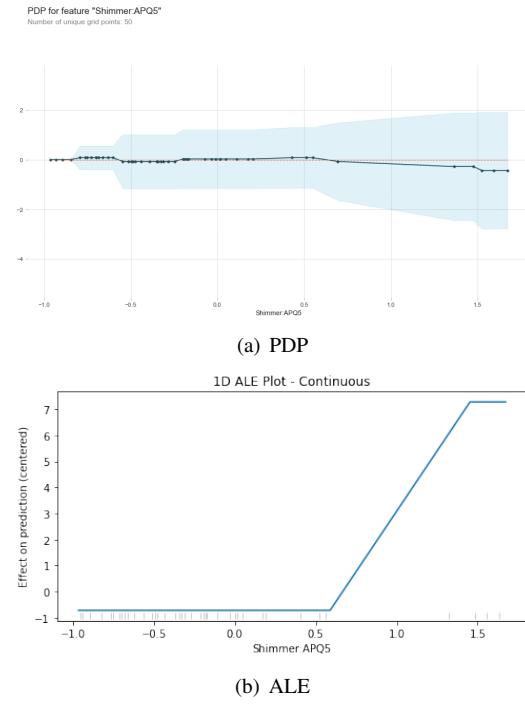


Fig. 166. PDP and ALE for the "Shimmer:APQ5" Feature in the K-Means based classifier

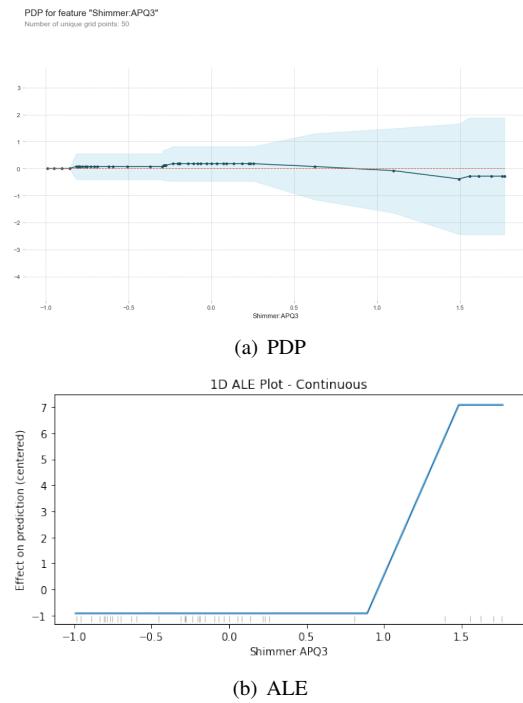


Fig. 165. PDP and ALE for the "Shimmer:APQ3" Feature in the K-Means based classifier

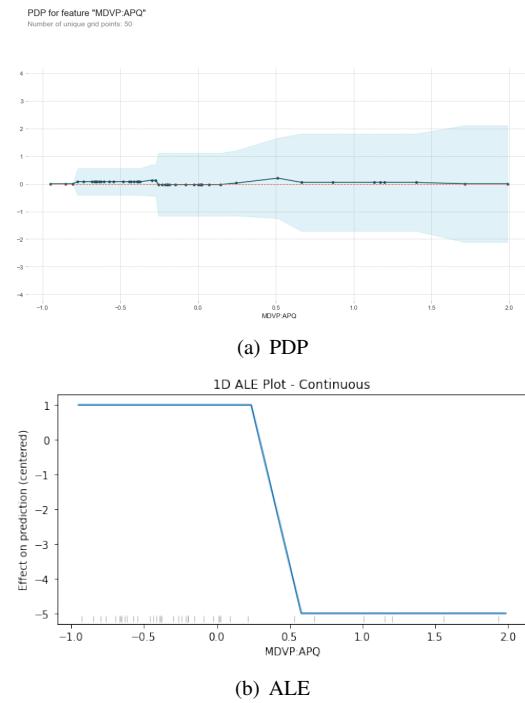
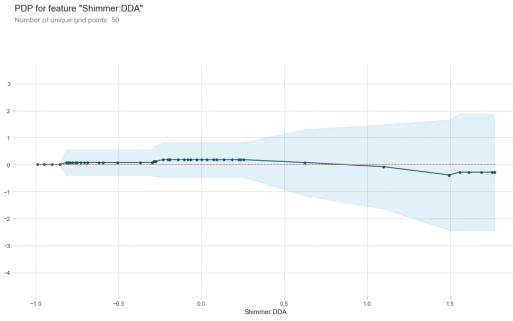
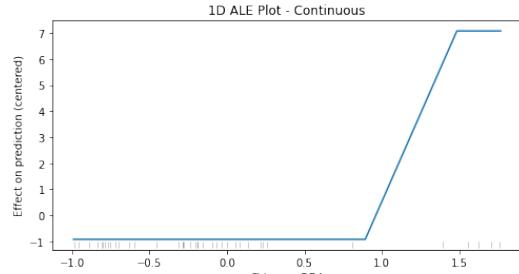


Fig. 167. PDP and ALE for the "MDVP:APQ" Feature in the K-Means based classifier

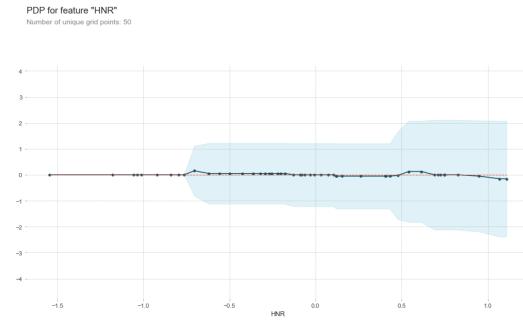


(a) PDP

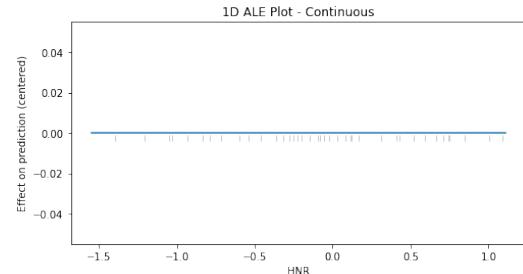


(b) ALE

Fig. 168. PDP and ALE for the "Shimmer:DDA" Feature in the K-Means based classifier

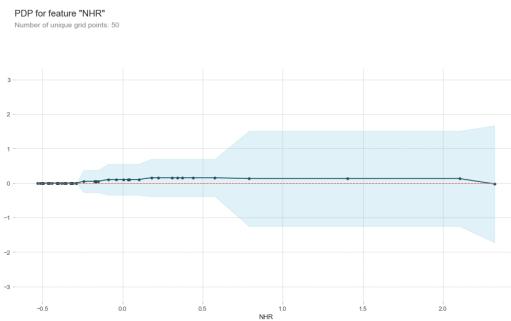


(a) PDP

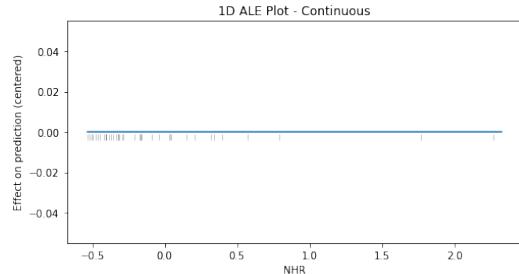


(b) ALE

Fig. 170. PDP and ALE for the "HNR" Feature in the K-Means based classifier

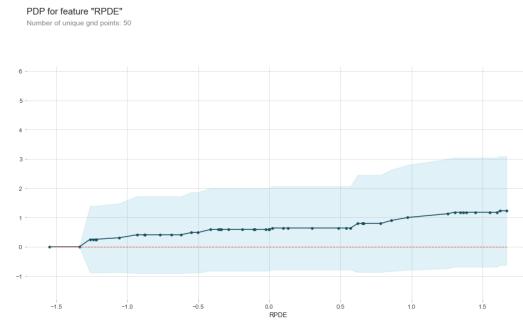


(a) PDP

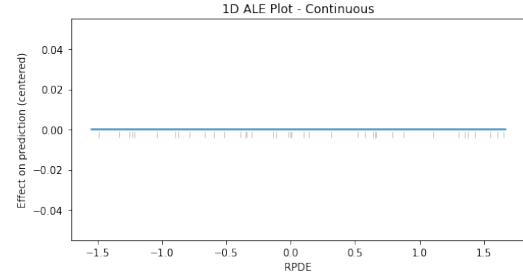


(b) ALE

Fig. 169. PDP and ALE for the "NHR" Feature in the K-Means based classifier

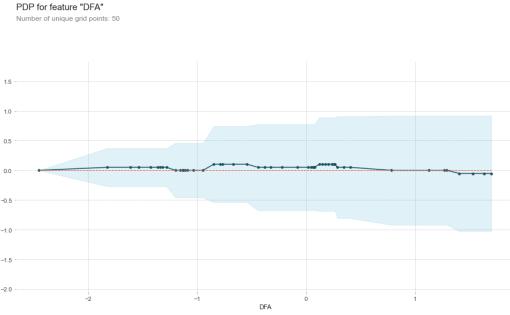


(a) PDP

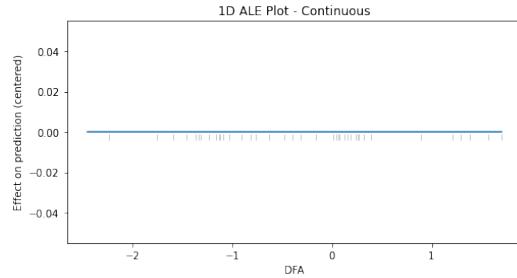


(b) ALE

Fig. 171. PDP and ALE for the "RPDE" Feature in the K-Means based classifier

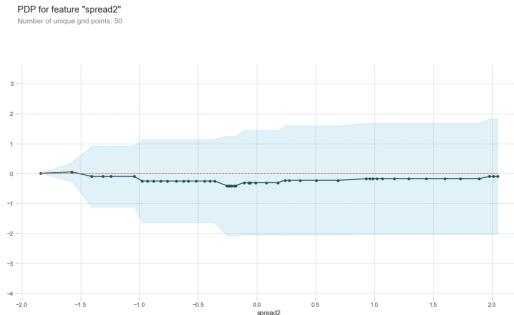


(a) PDP

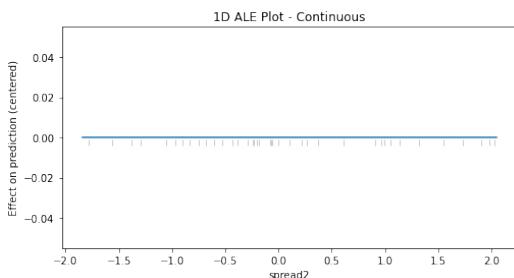


(b) ALE

Fig. 172. PDP and ALE for the "DFA" Feature in the K-Means based classifier

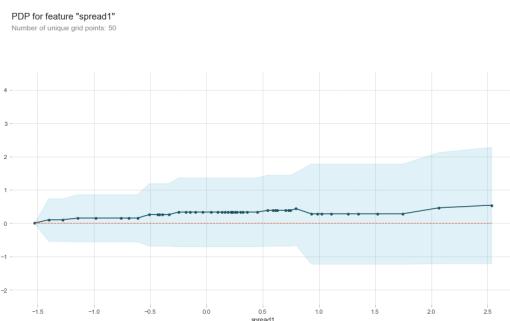


(a) PDP

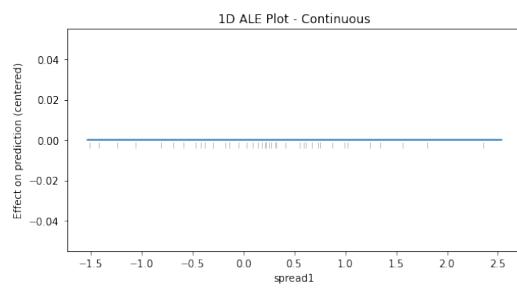


(b) ALE

Fig. 174. PDP and ALE for the "spread2" Feature in the K-Means based classifier



(a) PDP

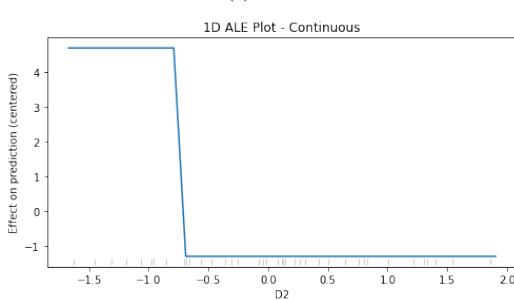


(b) ALE

Fig. 173. PDP and ALE for the "spread1" Feature in the K-Means based classifier



(a) PDP



(b) ALE

Fig. 175. PDP and ALE for the "D2" Feature in the K-Means based classifier

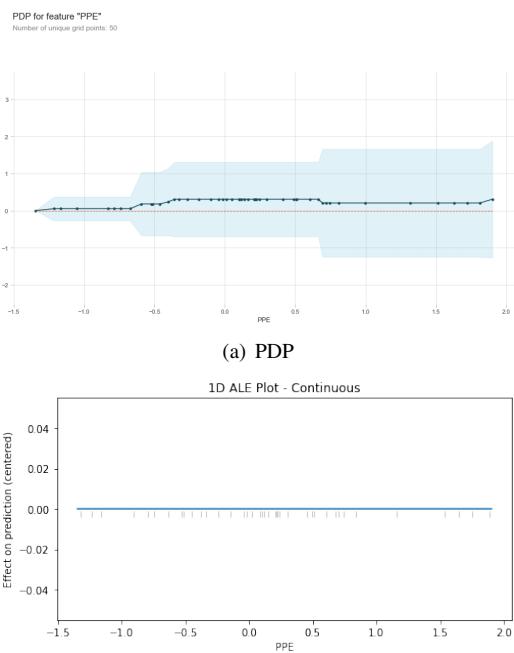
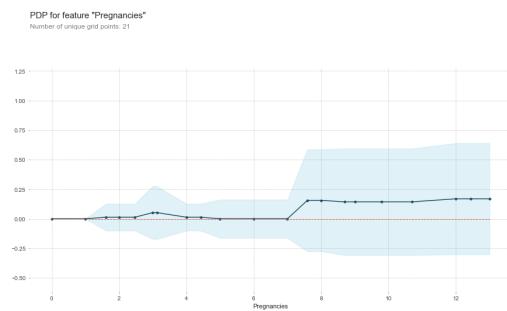
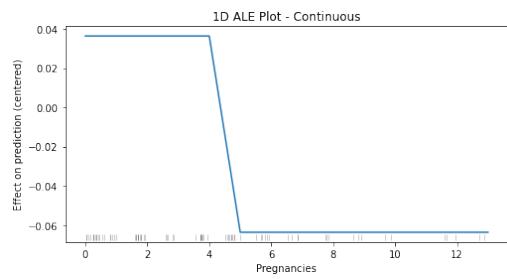


Fig. 176. PDP and ALE for the "PPE" Feature in the K-Means based classifier

S. ALE and PDP plots for the Decision Tree on the Diabetes

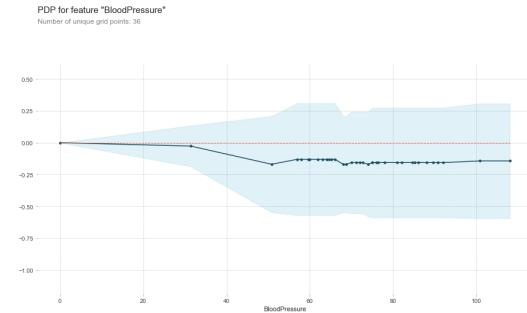


(a) PDP

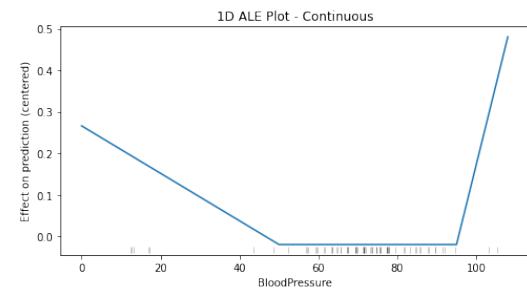


(b) ALE

Fig. 177. PDP and ALE for the "Pregnancies" Feature in the Decision Tree

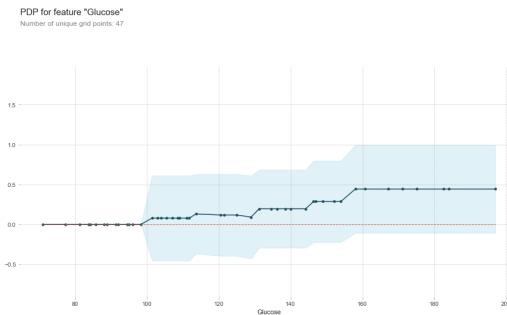


(a) PDP

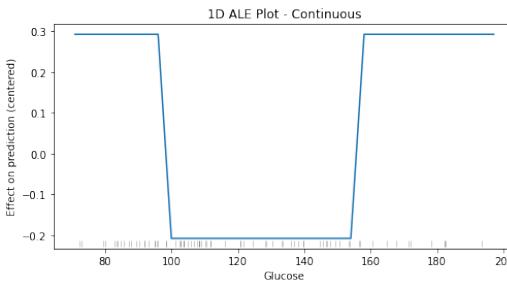


(b) ALE

Fig. 179. PDP and ALE for the "BloodPressure" Feature in the Decision Tree

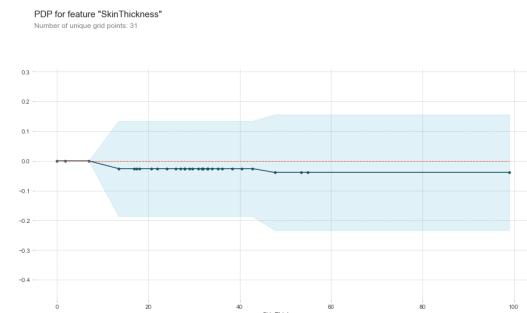


(a) PDP

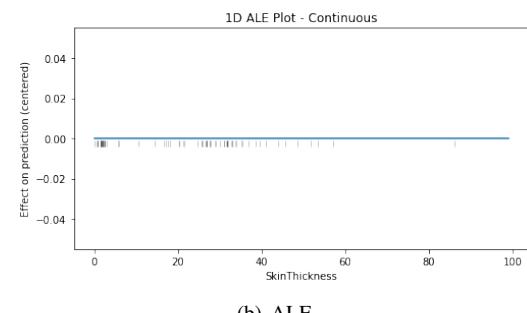


(b) ALE

Fig. 178. PDP and ALE for the "Glucose" Feature in the Decision Tree

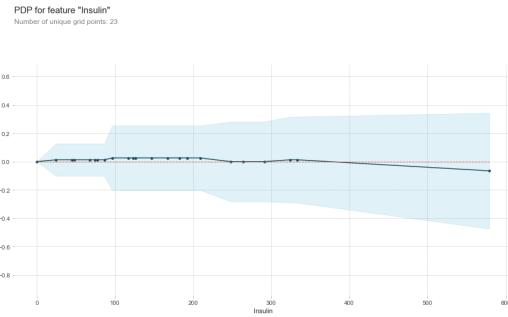


(a) PDP

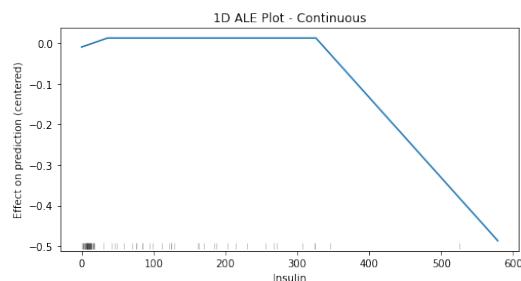


(b) ALE

Fig. 180. PDP and ALE for the "SkinThickness" Feature in the Decision Tree

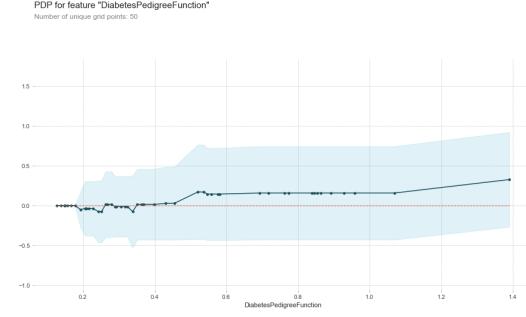


(a) PDP

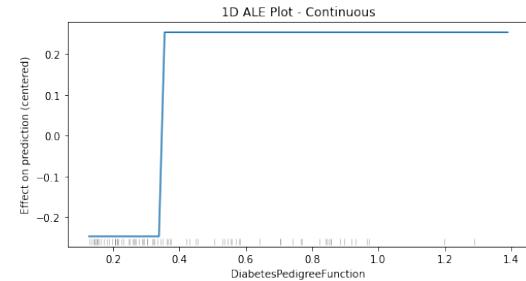


(b) ALE

Fig. 181. PDP and ALE for the "Insulin" Feature in the Decision Tree

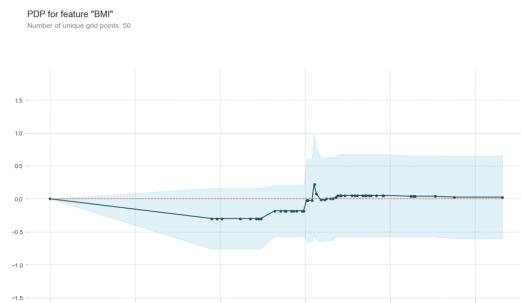


(a) PDP

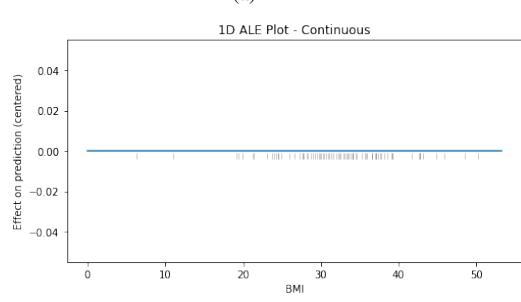


(b) ALE

Fig. 183. PDP and ALE for the "DiabetesPedigreeFunction" Feature in the Decision Tree



(a) PDP

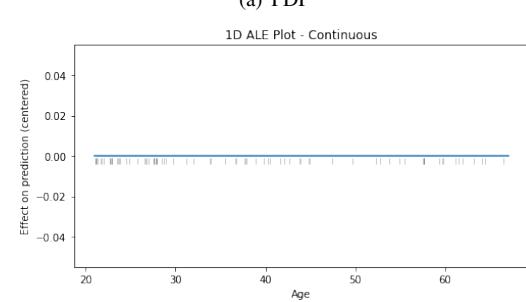


(b) ALE

Fig. 182. PDP and ALE for the "BMI" Feature in the Decision Tree



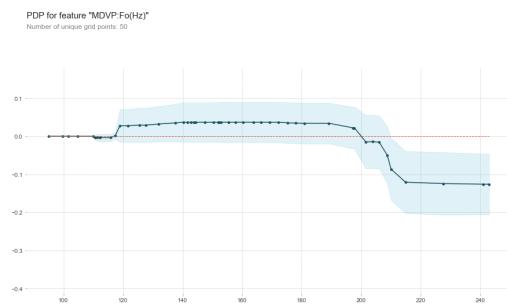
(a) PDP



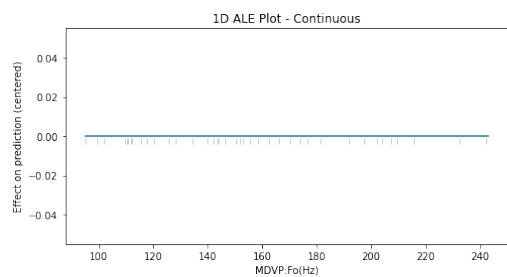
(b) ALE

Fig. 184. PDP and ALE for the "Age" Feature in the Decision Tree

T. ALE and PDP plots for the Random Forest on the Diabetes

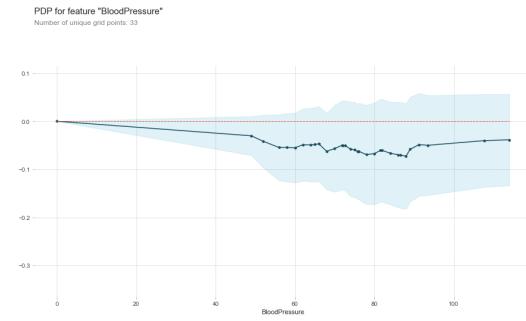


(a) PDP

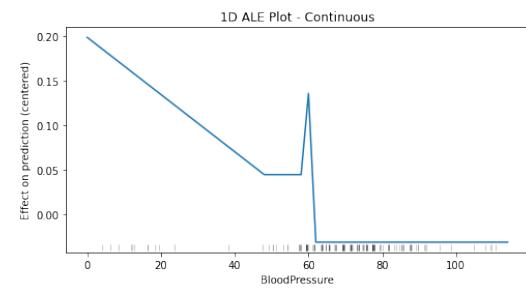


(b) ALE

Fig. 185. PDP and ALE for the "Pregnancies" Feature in the Random Forest

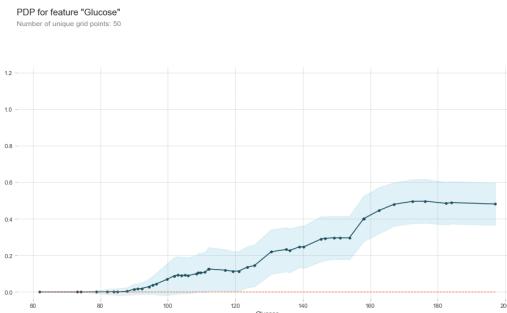


(a) PDP

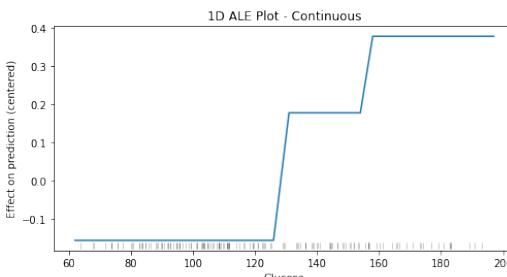


(b) ALE

Fig. 187. PDP and ALE for the "BloodPressure" Feature in the Random Forest

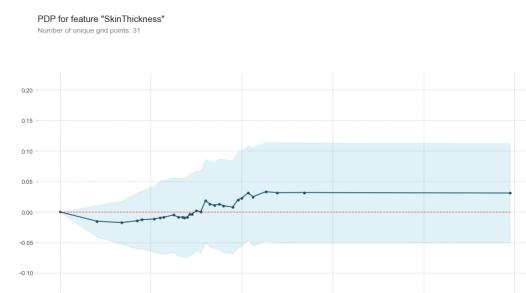


(a) PDP

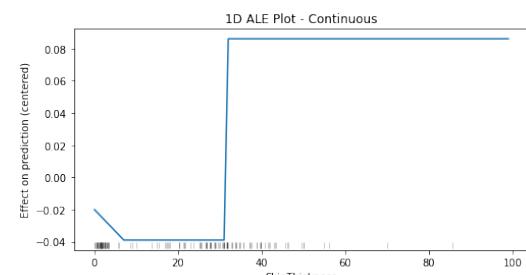


(b) ALE

Fig. 186. PDP and ALE for the "Glucose" Feature in the Random Forest



(a) PDP



(b) ALE

Fig. 188. PDP and ALE for the "SkinThickness" Feature in the Random Forest

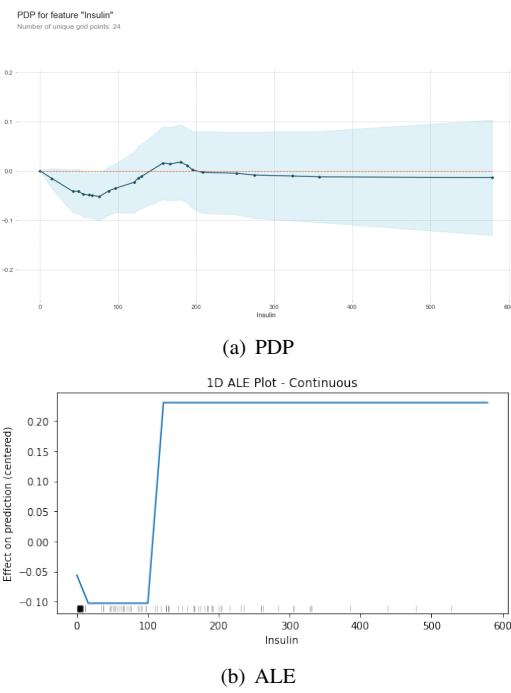


Fig. 189. PDP and ALE for the "Insulin" Feature in the Random Forest

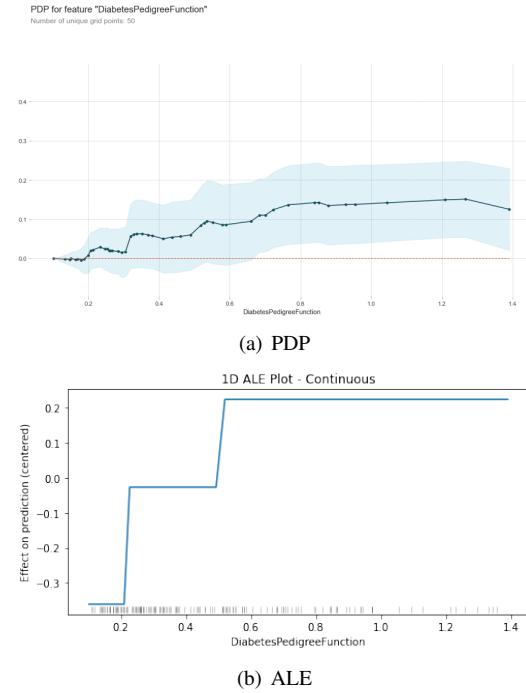


Fig. 191. PDP and ALE for the "DiabetesPedigreeFunction" Feature in the Random Forest

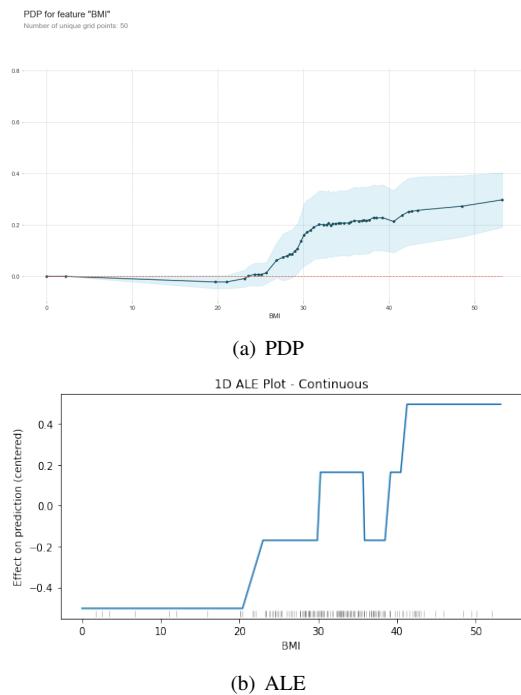


Fig. 190. PDP and ALE for the "BMI" Feature in the Random Forest

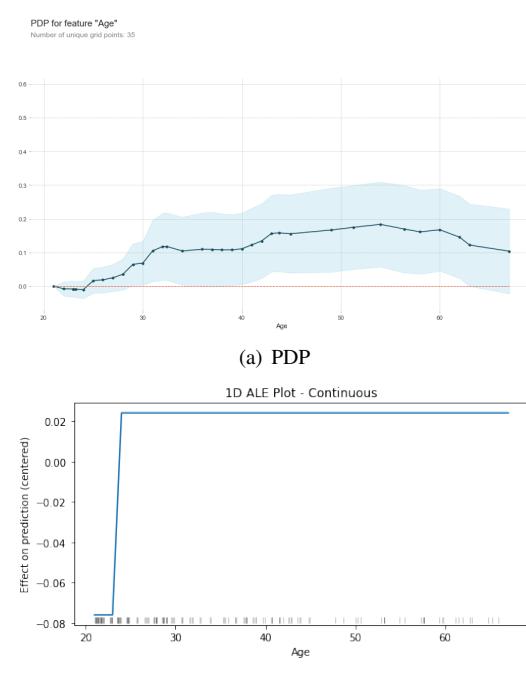


Fig. 192. PDP and ALE for the "Age" Feature in the Random Forest

U. ALE and PDP plots for the KNN on the Diabetes

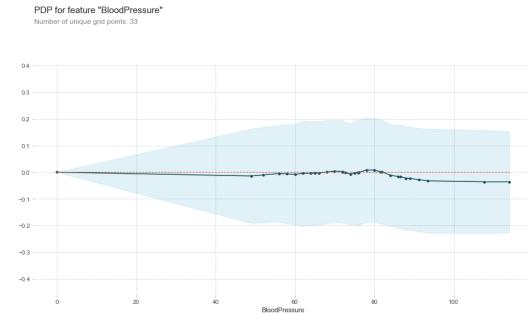
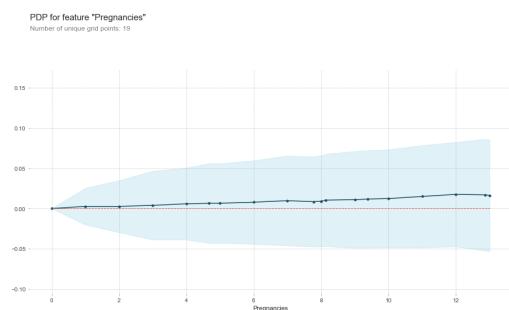


Fig. 193. PDP and ALE for the "Pregnancies" Feature in the knn

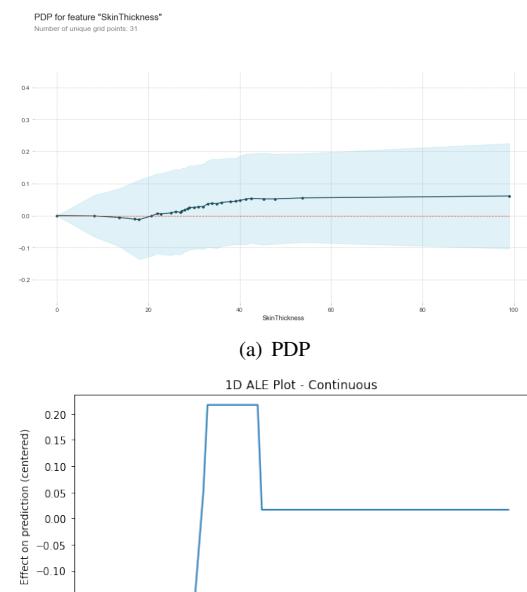
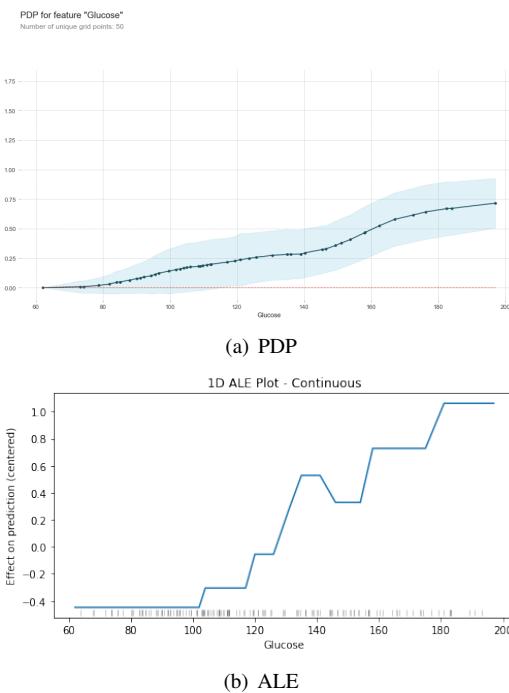
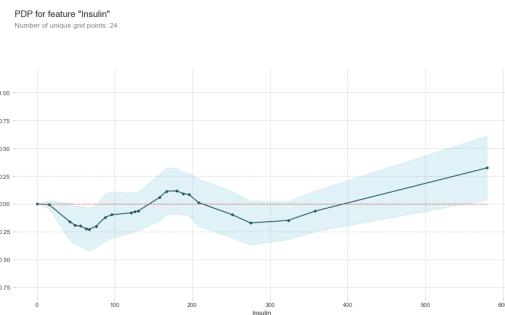


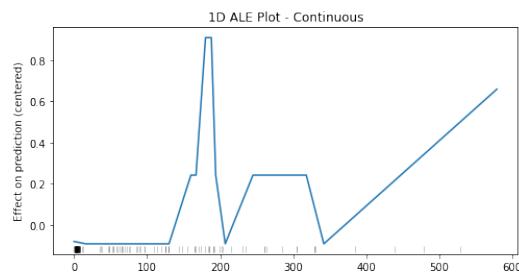
Fig. 194. PDP and ALE for the "Glucose" Feature in the knn

Fig. 195. PDP and ALE for the "BloodPressure" Feature in the knn

Fig. 196. PDP and ALE for the "SkinThickness" Feature in the knn

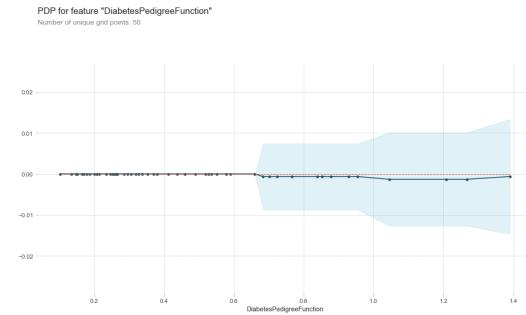


(a) PDP

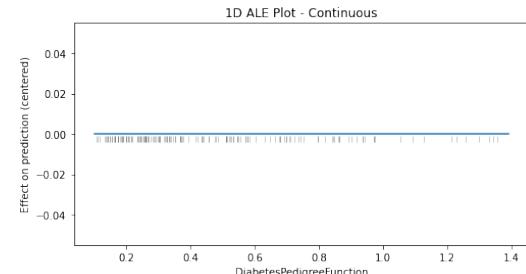


(b) ALE

Fig. 197. PDP and ALE for the "Insulin" Feature in the knn

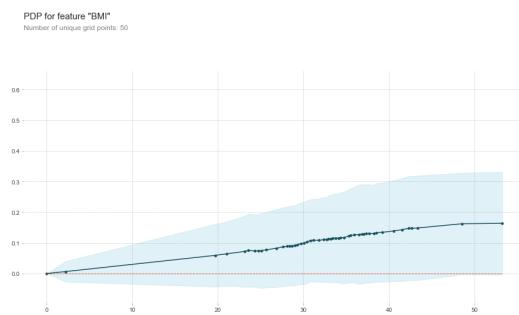


(a) PDP

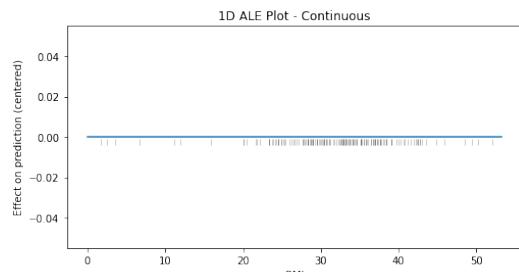


(b) ALE

Fig. 199. PDP and ALE for the "DiabetesPedigreeFunction" Feature in the knn

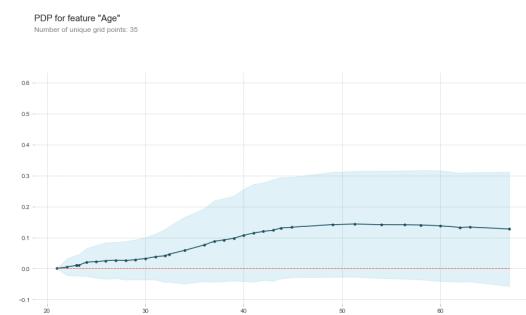


(a) PDP

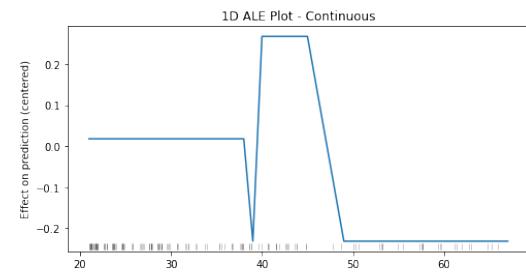


(b) ALE

Fig. 198. PDP and ALE for the "BMI" Feature in the knn



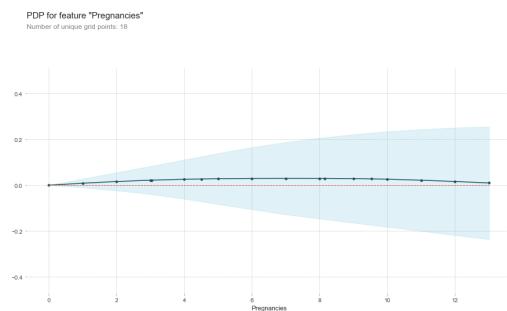
(a) PDP



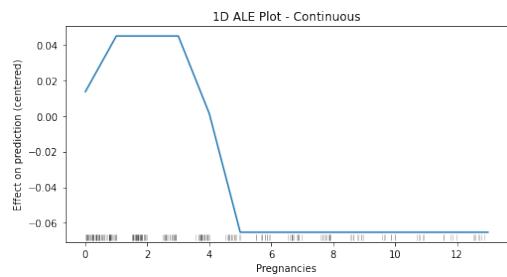
(b) ALE

Fig. 200. PDP and ALE for the "Age" Feature in the knn

V. ALE and PDP plots for the MLP on the Diabetes

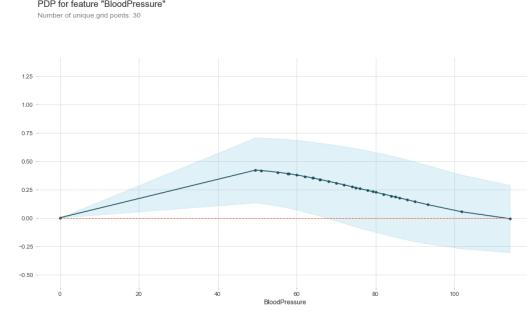


(a) PDP

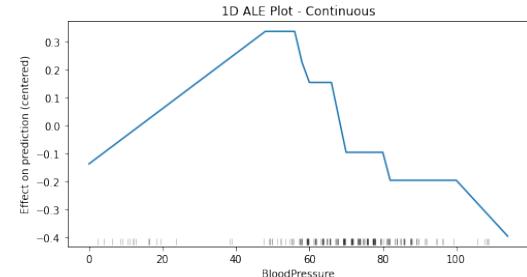


(b) ALE

Fig. 201. PDP and ALE for the "Pregnancies" Feature in the MLP

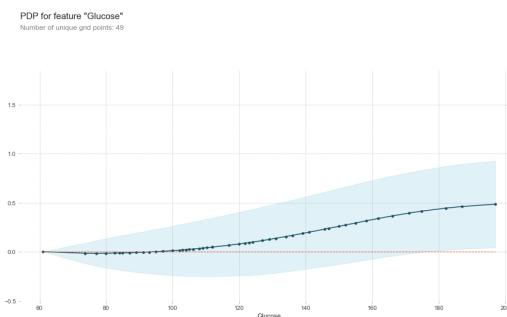


(a) PDP

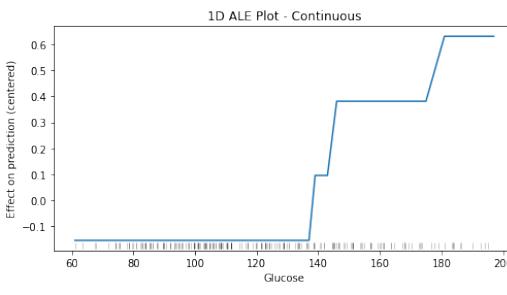


(b) ALE

Fig. 203. PDP and ALE for the "BloodPressure" Feature in the MLP

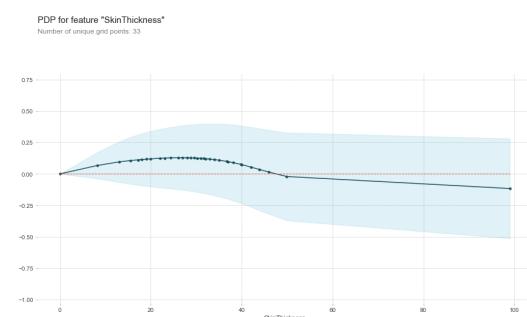


(a) PDP

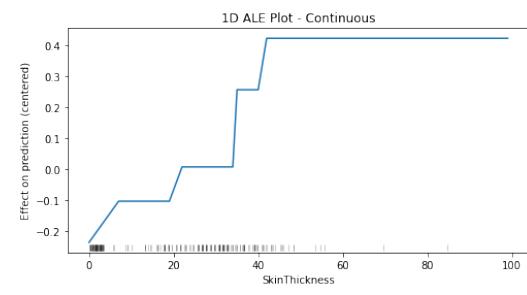


(b) ALE

Fig. 202. PDP and ALE for the "Glucose" Feature in the MLP

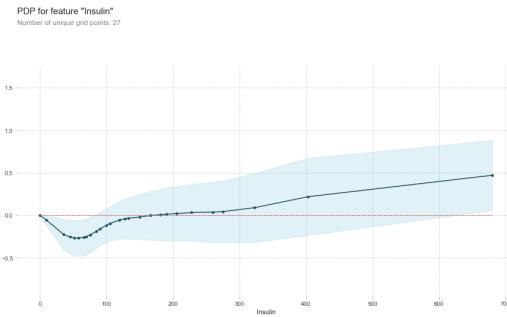


(a) PDP

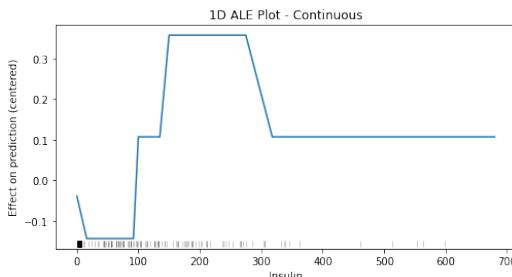


(b) ALE

Fig. 204. PDP and ALE for the "SkinThickness" Feature in the MLP

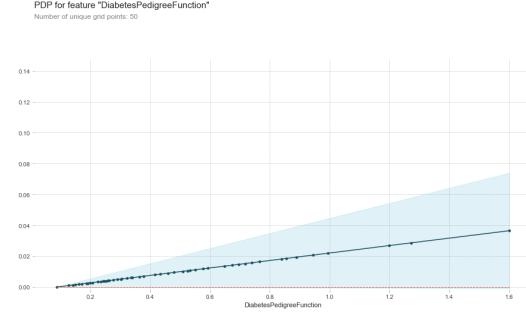


(a) PDP

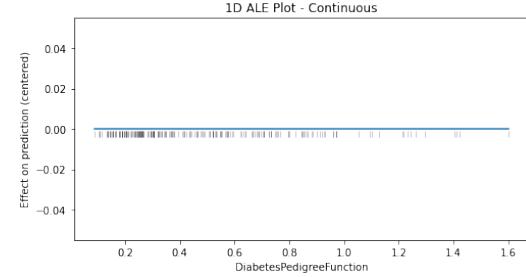


(b) ALE

Fig. 205. PDP and ALE for the "Insulin" Feature in the MLP

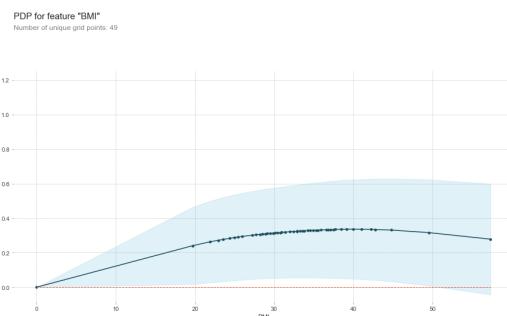


(a) PDP

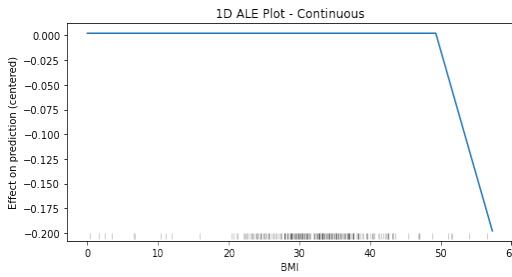


(b) ALE

Fig. 207. PDP and ALE for the "DiabetesPedigreeFunction" Feature in the MLP

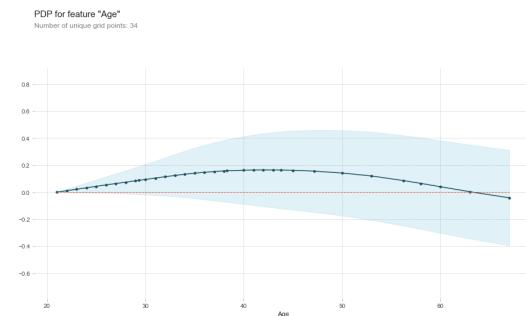


(a) PDP

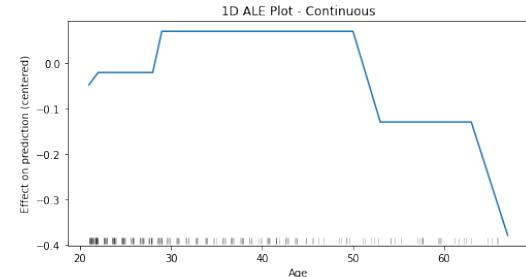


(b) ALE

Fig. 206. PDP and ALE for the "BMI" Feature in the MLP



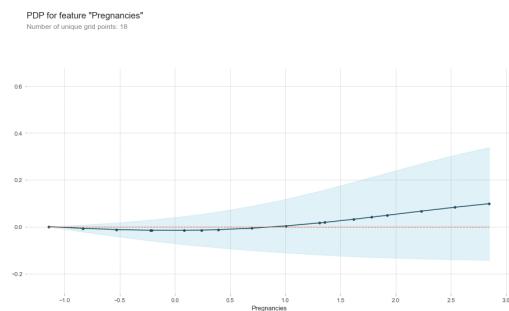
(a) PDP



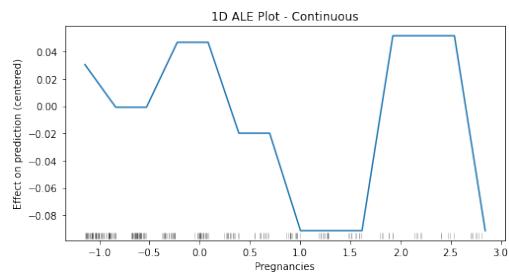
(b) ALE

Fig. 208. PDP and ALE for the "Age" Feature in the MLP

W. ALE and PDP plots for the SVM on the Diabetes

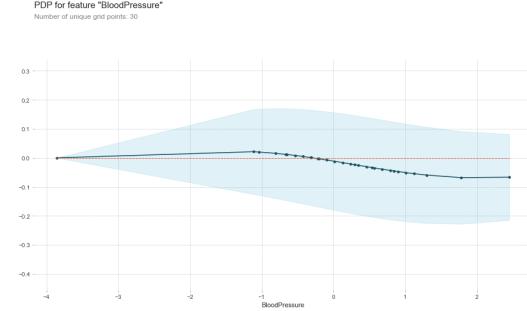


(a) PDP

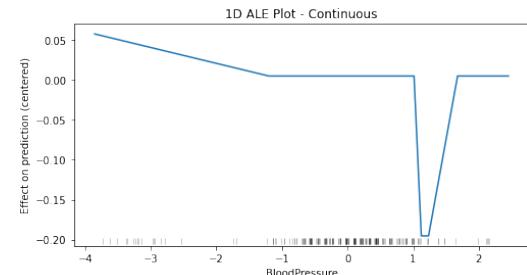


(b) ALE

Fig. 209. PDP and ALE for the "Pregnancies" Feature in the SVM

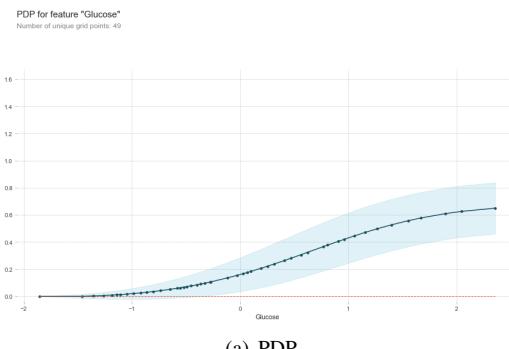


(a) PDP

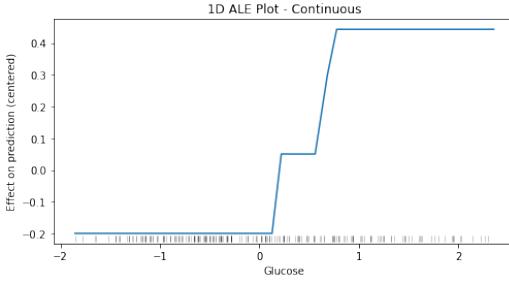


(b) ALE

Fig. 211. PDP and ALE for the "BloodPressure" Feature in the SVM

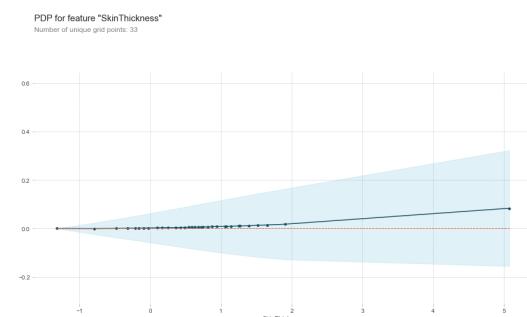


(a) PDP

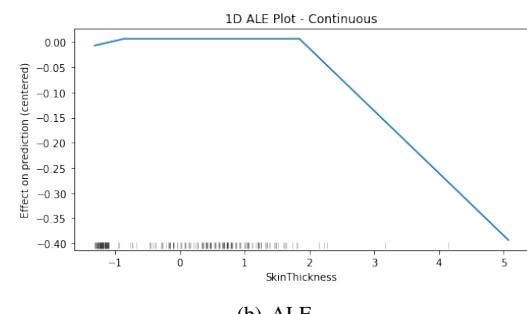


(b) ALE

Fig. 210. PDP and ALE for the "Glucose" Feature in the SVM

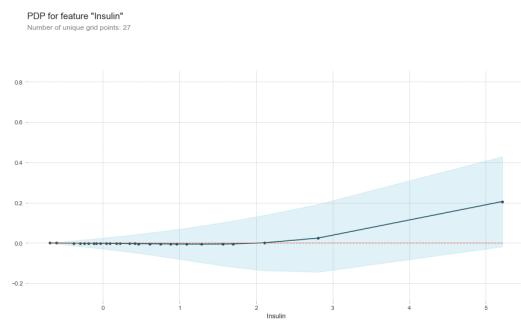


(a) PDP

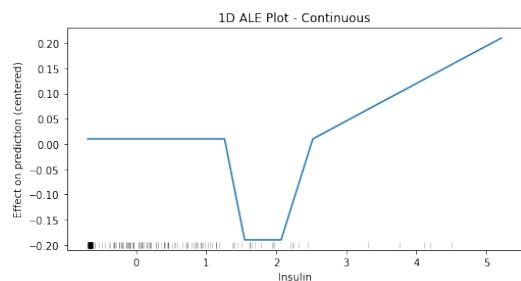


(b) ALE

Fig. 212. PDP and ALE for the "SkinThickness" Feature in the SVM

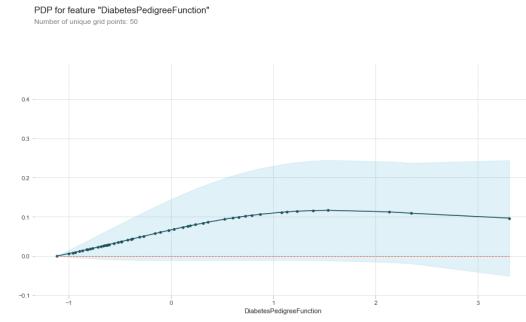


(a) PDP

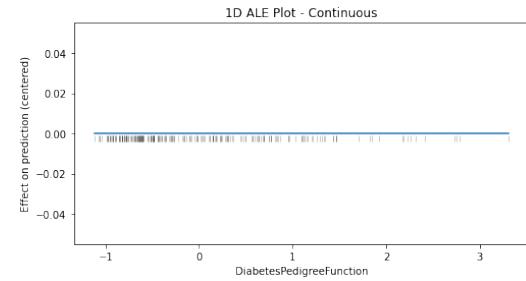


(b) ALE

Fig. 213. PDP and ALE for the "Insulin" Feature in the SVM

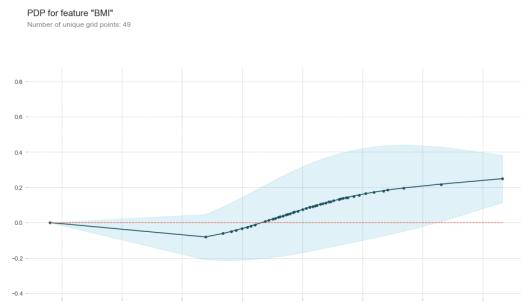


(a) PDP

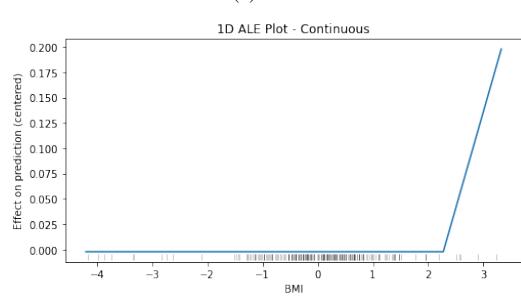


(b) ALE

Fig. 215. PDP and ALE for the "DiabetesPedigreeFunction" Feature in the SVM

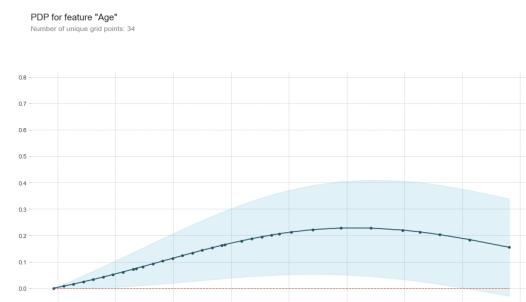


(a) PDP

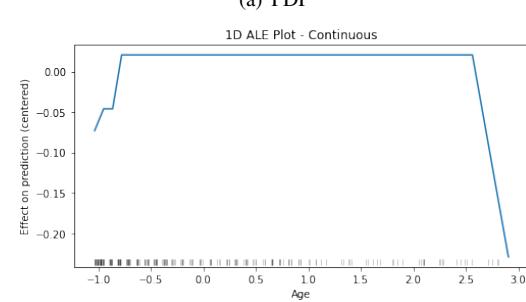


(b) ALE

Fig. 214. PDP and ALE for the "BMI" Feature in the SVM



(a) PDP



(b) ALE

Fig. 216. PDP and ALE for the "Age" Feature in the SVM

X. ALE and PDP plots for the K-Means based classifier on the Diabetes

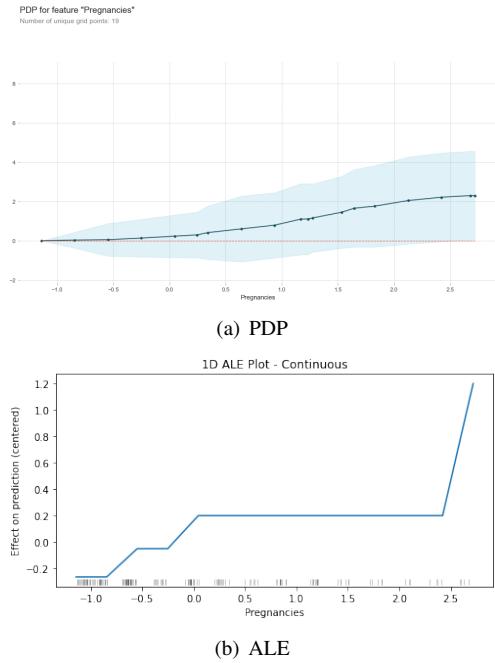


Fig. 217. PDP and ALE for the "Pregnancies" Feature in the K-Means based classifier

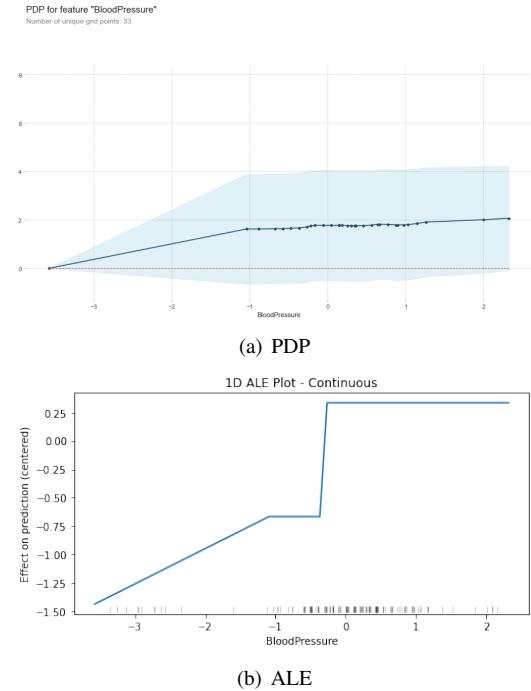


Fig. 219. PDP and ALE for the "BloodPressure" Feature in the K-Means based classifier

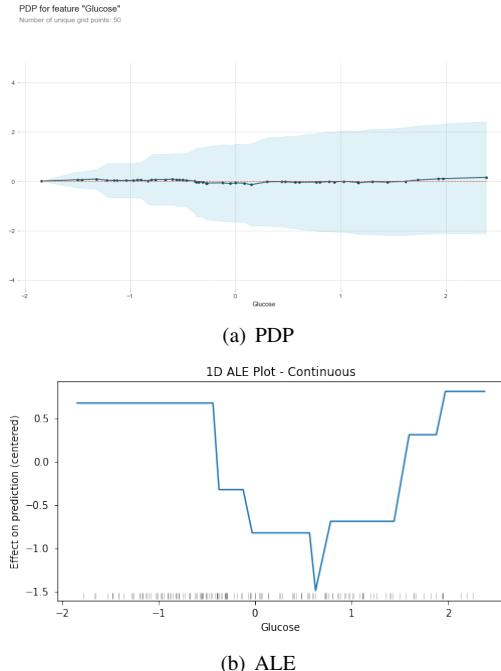


Fig. 218. PDP and ALE for the "Glucose" Feature in the K-Means based classifier

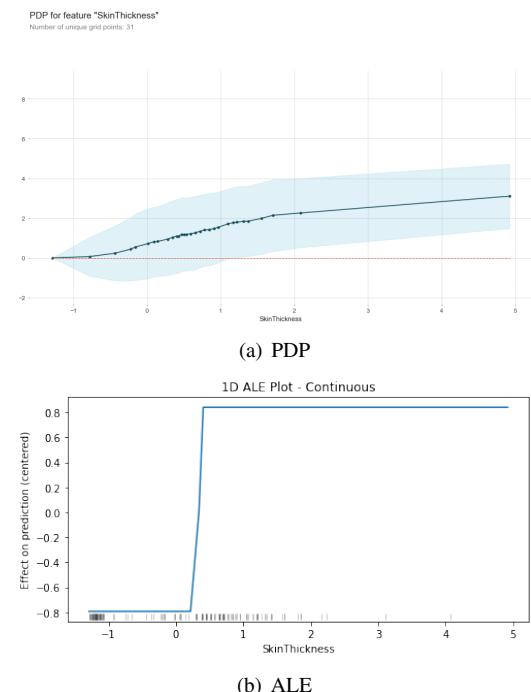
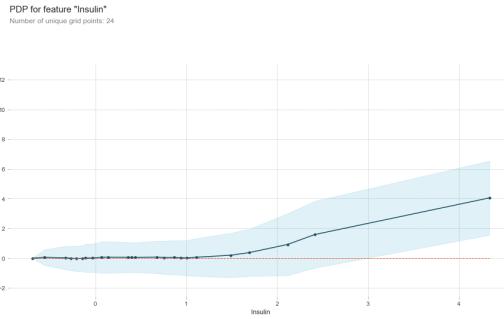
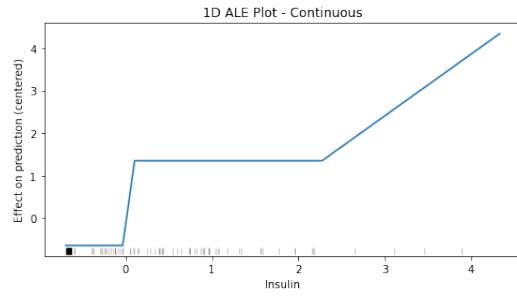


Fig. 220. PDP and ALE for the "SkinThickness" Feature in the K-Means based classifier

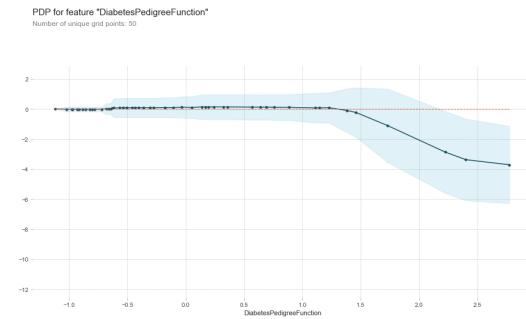


(a) PDP

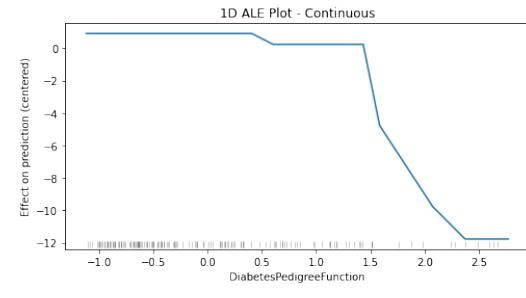


(b) ALE

Fig. 221. PDP and ALE for the "Insulin" Feature in the K-Means based classifier

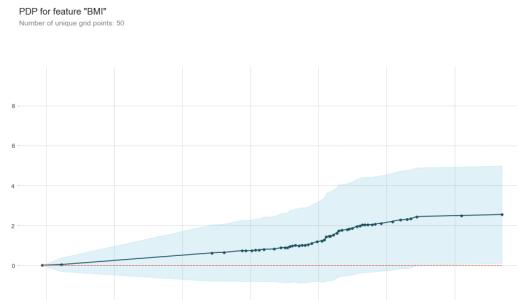


(a) PDP

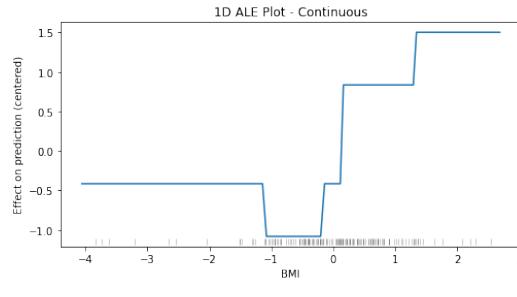


(b) ALE

Fig. 223. PDP and ALE for the "DiabetesPedigreeFunction" Feature in the K-Means based classifier

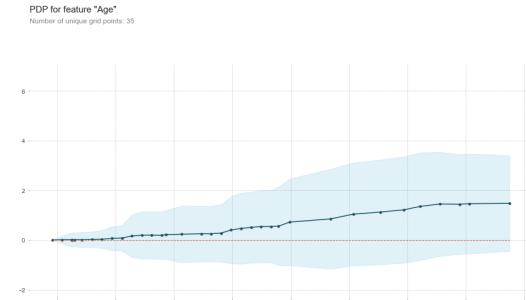


(a) PDP

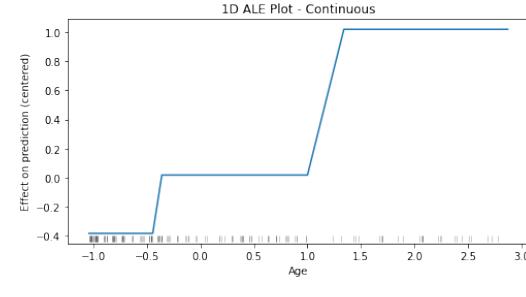


(b) ALE

Fig. 222. PDP and ALE for the "BMI" Feature in the K-Means based classifier



(a) PDP



(b) ALE

Fig. 224. PDP and ALE for the "Age" Feature in the K-Means based classifier