# Mixture of Decision Trees for Interpretable Machine Learning

Simeon Brüggenjürgen[1], Nina Schaaf[2], Pascal Kerschke[3], and Marco F. Huber[2,4]

*Abstract*—This work introduces a novel interpretable machine learning method called Mixture of Decision Trees (MoDT). It constitutes a special case of the Mixture of Experts ensemble architecture, which utilizes a linear model as gating function and decision trees as experts. Our proposed method is ideally suited for problems that cannot be satisfactorily learned by a single decision tree, but which can alternatively be divided into subproblems. Each subproblem can then be learned well from a single decision tree. Therefore, MoDT can be considered as a method that improves performance while maintaining interpretability by making each of its decisions understandable and traceable to humans.

Our work is accompanied by a Python implementation, which uses an interpretable gating function, a fast learning algorithm, and a direct interface to fine-tuned interpretable visualization methods. The experiments confirm that the implementation works and, more importantly, show the superiority of our approach compared to single decision trees and random forests of similar complexity.

## I. INTRODUCTION

Machine learning (ML) has been and continues to be an enormous success story. Historically, it has been the goal of researchers to develop problem-solving ML algorithms with ever-increasing accuracy. However, in recent years, the public and research interest in interpretable machine learning (IML) and its overarching field explainable artificial intelligence (XAI) has soared [1], [2]. This development is fueled by the desire of humans to understand how algorithms and machines operate in an increasingly digitized and data-driven world. While algorithms are developed that continuously claim more domains where they are outperforming humans, the inner workings of these algorithms tend to become more complex. As humans are arguably not even able to focus on more than one task at once, it becomes clear that it is almost impossible to comprehensively understand the plethora of internal factors that contribute to a complex model's decision. Still, humans want to use, evaluate, and trust algorithms. To accommodate the two latter desires, methods of IML can be used to simplify, distill, and translate algorithmic decisions into a human-comprehensible format.

There is a plethora of methods in the field of IML that can be categorized along multiple dimensions. A typical criterion is whether an IML method can be used on any underlying ML model (model-agnostic), or if it is only suitable for a specific type of model (model-specific) [3]. On another dimension, IML methods that are used after the creation of the to-be analyzed ML model are grouped under the term post-hoc interpretability [3]–[5]. Examples of model-agnostic post-hoc methods are LIME [6] and SHAP [7]. In contrast, there are IML methods where interpretability is already gained at the time of the creation of the ML model. These intrinsically interpretable models are by definition always model-specific [1], [8]. Examples of an intrinsically interpretable method are a (simple) linear regression, decision trees (DTs), or the here presented new method.

A DT is an ML technique that can be highly interpretable. This is due to the simple and intrinsic way to visualize a DT model and the ability to decompose it into easily understandable decision rules. A DT can be used to solve ML problems directly, however, it is not the most powerful method. Depending on the problem, there can be alternatives like deep neural networks or random forests that achieve higher performance but are in turn not interpretable anymore [8]. Therefore, research aims at developing algorithms that exceed the performance limitations of DTs but still remain interpretable. An approach by Vasic et al. [9] successfully merges DTs and the Mixture of Experts (MoE) architecture. MoE is an ensemble technique that was initially proposed in 1991 by Jacobs et al. [10]. It uses a so-called gating function to assign the input data to distinct experts. While there have been plenty of publications that explore the possibilities of MoE, the technique remains a niche. Vasic et al. [9] demonstrate that it is possible to successfully combine the technique with DTs, however, do not discuss the implications for IML.

In this paper we introduce a variant of the MoE architecture using DTs that specifically focuses on interpretability. To the best of our knowledge, such an approach has not been published before. A ready-to-be-used implementation of this approach is written in Python and can be found on GitHub[1]. In order to highlight the methods' reliance on DTs, the approach is given the name *Mixture of Decision Trees (MoDT)*. Our method builds on top of the one by Vasic et al. [9] who use a similar architecture, however, do not utilize its potential for interpretability. The main contributions of MoDT are a novel interpretable gating function that can reduce a multi-

[1]Simeon Brüggenjürgen is with the Department of Information Systems, University of Münster, Germany.
simeon.brueggenjuergen@uni-muenster.de

[2]Nina Schaaf and Marco F. Huber are with the Department Cyber Cognitive Intelligence (CCI), Fraunhofer Institute for Manufacturing Engineering and Automation IPA, Stuttgart, Germany. marco.huber@ieee.org

[3]Pascal Kerschke is with the Institute of Transport and Economics, Dresden University of Technology, Germany.

[4]Marco F. Huber is with the Institute of Industrial Manufacturing and Management IFF, University of Stuttgart, Germany.

[1]https://github.com/simsal0r/mixture-of-decision-trees

dimensional gating problem to a two-dimensional one, the usage of a fast linear regression for finding promising model parameters, and a visualization method that matches DTs and gating function. Thereby, MoDT is the first modification of the MoE architecture using DTs that is fully optimized towards interpretability.

In the following, fundamental techniques and a problem setting for MoDT are introduced. The developed method is described in detail in Sec. III. Next, MoDT is put onto the test bench with an experimental study that investigates its general performance boundaries. In Sec. V, the implications of this experiment and the usage of MoDT as a method of IML are discussed. Finally, the paper is briefly summarized.

## II. PROBLEM FORMULATION

In the following, it is assumed that a training dataset $\mathcal{D}$ is given comprising pairs $(x_i, y_i)$, with $i = 1 \ldots n$, where $x_i$ is the feature vector (extended by an element with value one) of dimension $p + 1$ and $y_i$ is the output (class label). All feature vectors are stored in a matrix $\mathbf{X} = [x_1 \ x_2 \ldots x_n]^{\mathrm{T}}$ and $y = [y_1 \ y_2 \ldots y_n]^{\mathrm{T}}$ is the target vector comprising all outputs. Then, a supervised classification model tries to learn the relationship between $x$ and $y$ such that it is able to predict $y'$ for vector $x'$ of a potentially different dataset $\mathcal{D}'$.

A DT is a supervised ML method that decomposes a problem into a series of decisions. This decomposition resembles the human decision process. Therefore, DTs are considered to be more accessible to humans than other techniques of ML and often given as a prime example of an IML method [1], [11]. Yet, while small trees are considered to be very interpretable, DTs can quickly become hard to understand if too many nodes (decision splits) are included. Therefore, the complexity of a DT is typically limited when interpretability is desired. This can be achieved by limiting the maximum tree depth $d$.

In 1991, the original MoE model was introduced by Jacobs et al. [10]. The general idea behind MoE is that multiple smaller problems are easier to solve than a single larger one. MoE can be seen as an ensemble method that uses multiple submodels (called *experts*) to solve a supervised ML problem. The number of experts $e$ needs to be set in advance. Then, each expert is trained independently on a subset of the input space. This is beneficial if the subsets incorporate different input-output relationships such that the experts can learn these distinct patterns. In MoE, the assignment of the subsets to the experts is conducted by a so-called *gating function*, which itself is subject to optimization. The gating function can be any function that maps an input data point to a subset of the input space. It is optimized alternately with the experts in an iterative fashion using the Expectation-Maximization (EM) algorithm [11]. Thereby, a "chicken and egg" dilemma can be avoided as the optimization of the gates depends on the optimization of the experts and vice versa. After training, when an MoE model is used for prediction, the gating function decides which expert is used for the prediction of a data point. Here, it is possible to only use the expert with the highest gating value or, alternatively, multiple experts whose outputs are combined
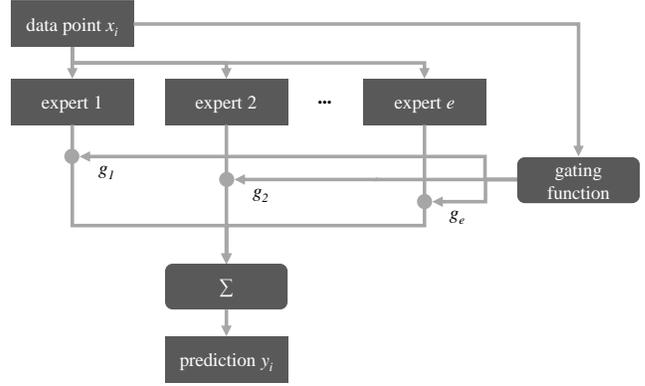


Figure 1. Architecture of an MoE prediction model for a single data point $x_i$. Depending on that point $x_i$, the gating function outputs weights which are used to aggregate the expert's individual predictions.

using weights provided by the gating function. Figure 1 depicts the architecture of a trained MoE model.

For an input data point $x_i$, the corresponding prediction $y_i$ of the MoE model is given by $y_i = \sum_{j=1}^{e} g_j \cdot z_j$ where $e$ denotes the number of experts, $g_j$ the gating weight for expert $j$, and $z_j$ the prediction of expert $j$ [11]. It holds that $g_j \in [0, 1]$ and $\sum_{j=1}^{e} g_j = 1$ [11]. This is often achieved through the usage of the softmax function. In the case where only one expert is used for the overall prediction, one weight is set to one while the others are set to zero. This is mostly a necessity for interpretability, as a combination of multiple experts can be difficult to understand for humans. Another necessity is that the gating function must be interpretable. Common MoE gating function choices like artificial neural networks are thus not appropriate for IML.

MoDT is an instance of MoE where a linear model is used as gating function and DTs are used as experts. A known problem of DTs is their inability to efficiently capture linear relationships due to the decision splits being always aligned with the dimensions of the input space. If two classes are optimally separated by a 45° boundary with respect to two dimensions, a DT can only approximate this with a large number of splits [11], which is not suitable in the context of interpretability. This problem can be addressed by MoDT by employing a linear gating function in conjunction with DTs.

## III. MIXTURE OF DECISION TREES

The key components of MoDT are introduced in the following. In Sec. III-A an overview of the entire MoDT architecture is given. Sec. III-B describes the training based on the EM algorithm, while Sec. III-C briefly addresses the prediction of the output $y$ given a trained MoDT and an input vector $x$. The visualization capabilities of MoDT to support interpretability are demonstrated in Sec. III-D by means of a toy dataset.

### A. Overview of the MoDT Architecture

The proposed method MoDT can be categorized as an ensemble method that uses multiple DTs as predictors. To

facilitate interpretability, only a small single-digit number of DTs should be used. Unlike many other ensemble methods like random forests (RFs) that combine the predictions of multiple DTs, MoDT only uses a *single* DT for a *single* prediction. Again, this is motivated by the requirement of interpretability: A single DT is usually considered to be interpretable while an aggregation of multiple DTs is not.

The number of regions and thus DTs $e$ has to be specified by the user. Alternatively, our implementation also includes methods to estimate a suitable number, e.g., based on Gaussian mixture models [11]. Each region is associated with a distinct DT that is trained on the subset of $\mathbf{X}$ that falls into the region. The core idea is that the regions contain distinct patterns that can be represented best by distinct DTs. To allow interpretability, the complexity of the DTs needs to be limited, which is achieved by setting a maximum DT depth $d$.

The gating parameter matrix $\theta \in \mathbb{R}^{(p+1) \times e}$ defines how the gating function splits the input space into regions and is initialized randomly. In MoDT, the gating function is a linear model including a bias term (intercept), which explains the additional dimension in $x_i$. The $j$-th column of $\theta$ comprises the parameters of the linear gating function corresponding to the $j$-th DT, with $j = 1, \ldots, e$. The outputs of the gating function, called *gating values* in the following, determine which parts of the input data are used for the training of each DT. The matrix of gating values $\mathbf{G} \in \mathbb{R}^{n \times e}$ is the result of a softmax on a matrix multiplication using the gating parameters $\theta$ and the input data $\mathbf{X}$ according to

$$\mathbf{G} = g(\mathbf{X}, \theta) = \begin{bmatrix} 1/|a_1| & 1/|a_1| & \cdots & 1/|a_1| \\ \vdots & \vdots & & \vdots \\ 1/|a_n| & 1/|a_n| & \cdots & 1/|a_n| \end{bmatrix} \circ \mathbf{A}, \quad (1)$$

where $\mathbf{A} = [a_1 \ a_2 \ \ldots \ a_n]^\mathrm{T} = \exp(\mathbf{X} \cdot \theta - \max(\mathbf{X} \cdot \theta))$, $a_i$ is the $i$-th row of $\mathbf{A}$, $\circ$ is the Hadamard (element-wise) product, and $g(.)$ is the *gating function*. Hence, each row of $\mathbf{G}$ sums to one. The calculation on the right-hand side of Eq. (III-A) corresponds to a row-wise application of a variation of the softmax function, where in each row of $\mathbf{X} \cdot \theta$ the row-wise maximum of $\mathbf{X} \cdot \theta$ is substracted inside the exponential function. The resulting gating values $\mathbf{G}$ determine the allocation of each data point to the regions. One element $g_{ij}$ of $\mathbf{G}$ represents the probability of choosing expert $j$ for data point $i$.

For MoDT, it is possible to use a gating function that either considers all $p$ features of $\mathbf{X}$, or alternatively, just two features. Due to their linear nature, both variants behave relatively predictable. Yet, the two-dimensional (2D) gating function is beneficial for interpretability as it has the great advantage of being comprehensibly visualizable using a 2D plot. For the 2D gating function, the question arises, which two out of all $p$ features should be selected for the gating function. In the implementation of MoDT, multiple feature importance methods (for example based on DTs, linear models, or principal component analysis) are included for selecting promising candidates. From an interpretability perspective, users can also manually select two features that are well known to them. The reduced dimension in case of a 2D gating function affects the size of $\theta$ and only the two selected feature columns of $\mathbf{X}$ are considerede for calculating $\mathbf{A}$ in (III-A). In our experiments (see Section IV), it will be tested how the limitation to two features impacts performance in comparison to the "full" gating function using all features.

The "goodness" of a region depends on the "goodness" of the associated DT and vice versa. To avoid a "chicken and egg" dilemma, the training algorithm thus uses the two-step EM algorithm, which alternately updates the gating function and the DTs. More precisely, during the E-step of the EM algorithm, a DT is trained for each region. Technically, this is achieved by using the output vector of the gating function for each data point $x$ as weight vector during the training of the DTs. Then, the expected "goodness" of the current combination of DT and the gating parameters $\theta$ with respect to the true prediction targets $y$ is calculated. This expectation is used in the M-step to optimize $\theta$ and the gating function, respectively. The E-step and M-step are repeated until a stopping criterion or a fixed number of iterations is reached.

### B. Training

In contrast to other MoE architectures, where the training of the experts is conducted in the M-step [11], [12], we prefer training the experts during the E-step. Further, regardless of the choice of the 2D or the full gating function, the following training procedure is the same.

During the training, the complexity of the DTs is controlled by a hyperparameter $d$ that limits the maximum depth of the DTs. There are also other possibilities to limit the complexity of a DT (like the number of leaf nodes), however here, the maximum depth is preferred as it potentially reduces the variance of the DTs between iterations, resulting in a more stable training process.

Once the DTs are trained, the expectation $E$, which gives the E-step its name, can be calculated. For this purpose, the gating function and a so-called *confidence function*, denoted as $c(\mathcal{T}, \mathbf{X}, y)$, are necessary, where $\mathcal{T}$ is the set of all DTs. At first, the gating values $\mathbf{G}$ are calculated by means of Eq. (III-A). Secondly, the confidence function $c(\mathcal{T}, \mathbf{X}, y)$ outputs the DTs' confidences for the correct predictions. For this purpose, as a first step, all trained DTs $z_j \in \mathcal{T}$ perform predictions for the input $\mathbf{X}$. Instead of assigning a particular class to a data point $x_i$, each DT $z_j$ calculates a vector comprising the prediction probabilities for all possible classes of $x_i$. This vector and the true classes $y$ are then used to extract the probability for a correct prediction of each DT for each data point of $\mathbf{X}$. Thus, the output of $c(.)$ corresponding to one data point $x_i$ is the vector $c_i$, where the $j$-th element of $c_i$ corresponds to the probability of DT $z_j$ outputting the true class $y_i$ of input vector $x_i$. Finally, the expectation $\mathbf{E}$ can be calculated according to

$$\mathbf{E} = \begin{bmatrix} 1/|b_1| & 1/|b_1| & \cdots & 1/|b_1| \\ \vdots & \vdots & & \vdots \\ 1/|b_n| & 1/|b_n| & \cdots & 1/|b_n| \end{bmatrix} \circ \mathbf{B}, \quad (2)$$

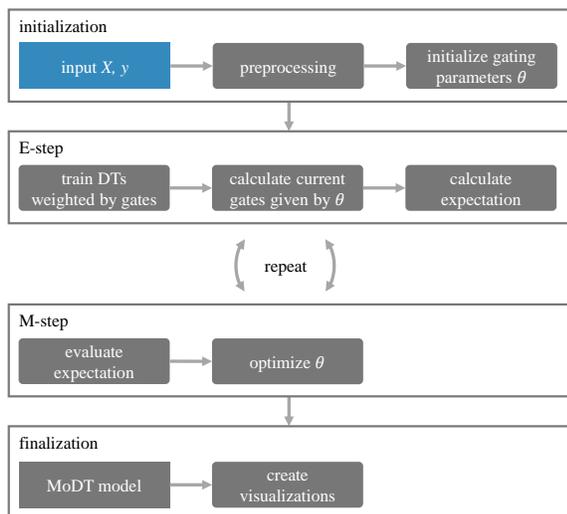Figure 2. Overview of the training process of MoDT.

where $\mathbf{B} = [b_1 \ b_2 \ \ldots \ b_n]^{\mathrm{T}} = \mathbf{G} \circ \mathbf{C}$ with $\mathbf{C} = [c_0 \ c_1 \ \ldots \ c_e] \in \mathbb{R}^{n \times e}$ being the output matrix of function $c(\mathcal{T}, \mathbf{X}, y)$. Each row of $\mathbf{E}$ sums to one.

The calculation of one row of $\mathbf{E}$ in Eq. (III-B) is illustrated in a numeric example: Given an MoDT model with three experts, assuming the gating values for a data point being $[0.7 \ 0.2 \ 0.1]$, and the experts' probabilities for the correct prediction are given by $[1.0 \ 0.9 \ 0.5]$, the expectation for this data point is $[0.70 \ 0.18 \ 0.05] \ / \ 0.93 = [0.753 \ 0.194 \ 0.054]$. If these values were directly interpreted as new gating values, it can be noticed that the probability for the best-performing expert increases, while the opposite is true for the weaker experts. This observation is crucial for the training algorithm of MoDT.

The newly calculated expectation does not replace the gating values but is used as a basis for a controllable optimization during the M-step of the EM algorithm. The original gating values and the expectation are used to calculate an optimization direction similar to a gradient that is applied to the gating parameters $\theta$, governed by an adjustable learning rate parameter. In classic MoE architectures, the expectation is used within a log-likelihood function that then serves as the loss for iterative and possibly computationally expensive optimization algorithms [11], [12]. However, an experiment by Yang and Ma [12] compares multiple optimization algorithms and concludes that a rather simple least squares linear regression is favorable for an MoE architecture with a linear gating function. Therefore, MoDT primarily uses a linear regression for the update of the gating parameters $\theta$ which does not require a gradient descent. The linear regression creates a linear model with coefficients $\beta$ that fits $\mathbf{X}$ to a target. For this purpose, the residual sum of squares (RSS), i.e., the sum of squared deviations of the predicted targets to the correct targets is minimized. The target of the regression is not the earlier calculated expectation $\mathbf{E}$, but the difference of $\mathbf{E}$ and the gating values $\mathbf{G}$. This target $\mathbf{E} - \mathbf{G} \in \mathbb{R}^{n \times e}$ can be

interpreted as the direction towards the desired new gating values. Formally, the coefficients $\beta$ are chosen to minimize

$$\mathrm{RSS}(\beta) = \sum_{i=1}^{n} \left( (e_i - g_i) - x_i^{\mathrm{T}} \beta \right) \cdot \left( (e_i - g_i) - x_i^{\mathrm{T}} \beta \right)^{\mathrm{T}} ,$$

with $e_i$ and $g_i$ being the $i$-th row of $\mathbf{E}$ and $\mathbf{G}$, respectively [13]. The resulting coefficients $\beta \in \mathbb{R}^{(p+1) \times e}$ can be seen as the internals of a linear model that matches $\mathbf{X}$ to $\mathbf{E} - \mathbf{G}$, similar to the linear gating function of MoDT which matches $\mathbf{X}$ to $\mathbf{G}$.

To calculate the new gating parameters $\theta_{\mathrm{new}}$, the update equation $\theta_{\mathrm{new}} = \theta_{\mathrm{old}} + \gamma \cdot \beta$ is employed, with $\theta_{\mathrm{old}}$ being the current gating parameters and $\gamma > 0$ being the learning rate. The update of $\theta$ concludes the M-step.

Once the EM algorithm has terminated, the trained MoDT model can be used to create visualizations for both the gating function (if the 2D variant is selected) and the DTs. An overview of the training process is given in Figure 2.

### C. Prediction

The prediction process takes unseen data $\mathbf{X}'$ as input of the gating function which outputs the gating values $\mathbf{G}$. It is important to note that the earlier outlined calculation of $\mathbf{G}$ (see Eq. (III-A)) is only used for the training and the gating function for the prediction is slightly different. In particular, the gating function for the prediction will only output the expert (DT) with the highest probability. The final prediction $y_i$ for a data point $x_i$ is thus given by

$$y_i = \sum_{j=1}^{e} z_j(x_i) \cdot \mathbb{1}(j = \arg \max_k (g_{ik})) ,$$

where $z_j$ is the prediction of the $j$-th DT and $\mathbb{1}$ is the indicator function. The indicator function returns one if a proposition is true, i.e., if data point $x_i$ belongs to expert $e_j$. Otherwise, zero is returned.

In contrast, it would also be possible to use the probabilities of the experts in place of the argmax function. In this case, the final prediction would be the result of a weighted sum of the predictions of all DTs. This architecture would resemble the common structure of MoE as seen earlier in Figure 1. This setup is also used in the paper of Vasic et al. [9] where it is called soft-gate prediction. However, for MoDT, such a combination of DTs is avoided as the increased complexity is unsuited for interpretability.

### D. Visualization and Application Example

In order to allow interpretability, a reasonable and comprehensive visualization is indispensable. If the feature space is two-dimensional or the 2D gating function is used, the gating function and the resulting regions can be plotted as depicted in Figure 3. Given the plot in Figure 3, a user can understand that MoDT solves the classification problem by dividing the dataset into three regions that incorporate three distinct patterns. A legend that enumerates the regions/DTs (Figure 3, top right) and a legend for the classes of the data points (Figure 3, bottom right) is added to the plot. The colors are chosen such that the data points remain visible on top of the regions. One
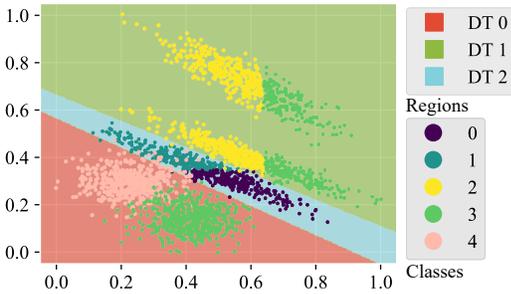
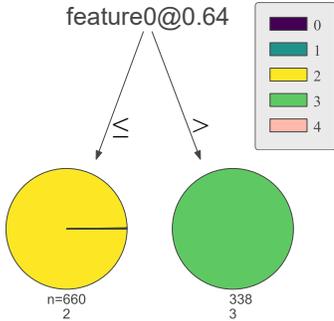Figure 3. Exemplary decision area of the 2D gating function.



Figure 4. DT corresponding to the green region in Figure 3.



Figure 5. Decision area of the 2D gating function on the steel dataset.



Figure 6. An example of a classic DT that is not very interpretable.

of the three DTs, i.e., the DT for the top-most green region can be seen in Figure 4. Here, the region is expressed by a simplistic DT with a single decision criterion, which often is called a decision stump. The top-right legend corresponds to the classes of the classification problem. The DTs for the blue and red regions are not depicted.

An example of the visualization of the 2D gating function on a more complex problem, namely the steel dataset [14], is depicted in Figure 5. The two features `Log_Y_Index` and `Steel_Plate_Thickness` have been selected using feature importance. Again, three experts are used. Although the corresponding DTs are not depicted in this paper, it can be seen how the regions incorporate distinct behavior. On this dataset, a MoDT model can reach a training accuracy of around $68\%$ when three DTs with a maximum depth of two are used (see Chapter IV-B). Here, MoDT uses at most $3 \cdot 7 = 21$ nodes (a DT of depth two has a maximum of seven nodes). For comparison, as a negative example, Figure 6 depicts a single DT that reaches the same accuracy. It uses almost twice the number of nodes (39) and—in contrast to MoDT—is too large to be conveniently plottable and interpretable. This dataset is revisited also in the experiments in Sec. IV.

For visualizing DTs, many software packages are freely available. Although the Python library dtreeviz (https://github.com/parrt/dtreeviz) does not check all interpretability boxes, it is a good match for MoDT. A great advantage of the library is that the colors of the classes can be adjusted with ease. This way, the colors of the DT correspond to the colors of the classes in the gating plot, thereby leveraging the joint comprehensibility of both graphics.
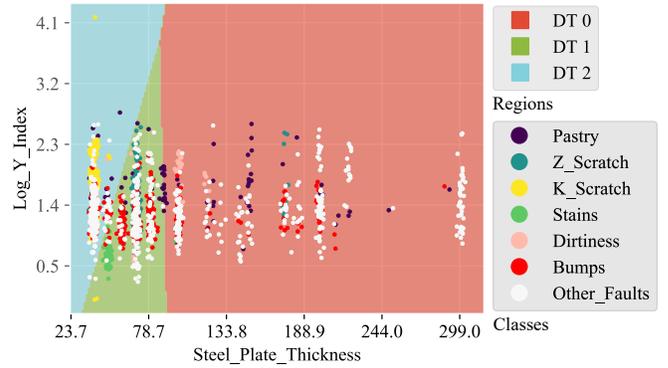
## IV. Experiments

In the following, the general performance of MoDT in terms of classification accuracy is tested and compared against itself and other classic tree-based ML models with varying complexity. MoDT is also compared to non-interpretable ML methods which generally should not be used for typical IML scenarios in which safety and ethical standards must be guaranteed [15]. However, as these methods are typically way more powerful regarding pure predictive performance, it can already be seen as a success if an interpretable method comes close to the performance of a non-interpretable method.

### A. Experimental Setup

To estimate an upper bound of MoDT's performance, a hyperparameter optimization algorithm is used that aims at finding a well-performing combination of hyperparameters. For this, the Python library Optuna [16] is utilized. The experiment is conducted with eight publicly available datasets [14], summarized in Table I. As some of the datasets contain categorical features that are not supported by the scikit-learn DT implementation, MoDT uses a one-hot-encoding for categorical features.

Each dataset is shuffled into training data (75%) and test data (25%). Then, the experiment is conducted independently for each dataset. Each training dataset is firstly used for the hyperparameter search which is repeated 500 times. Instead of choosing the single best hyperparameter combination, the ten

Table I
DATASETS USED IN THE EXPERIMENT.

| dataset | data points | classes | features (num./cat./enc.) | | | description |
|---|---|---|---|---|---|---|
| abalone | 4,177 | 3 | 7 | 1 | 10 | Predict age group of abalones (sea animals) based on physical measurements. |
| banknote | 1,372 | 2 | 4 | 0 | 4 | Predict authenticity of banknotes. |
| breast | 569 | 2 | 10 | 0 | 10 | Predict breast cancer based on properties extracted from images. |
| cars | 1,728 | 4 | 0 | 6 | 21 | Predict quality of cars based on categorical features only. |
| contraceptive | 1,473 | 3 | 2 | 7 | 24 | Predict contraceptive method. |
| iris | 150 | 3 | 4 | 0 | 4 | Predict species of a flower. Popular dataset. |
| steel | 1,941 | 7 | 27 | 0 | 27 | Predict defects of steel plates. |
| students | 666 | 4 | 0 | 11 | 49 | Predict performances of students on an entrance examination. |

best combinations are identified. These are then used to train ten distinct MoDT models on the training data. Lastly, the final performance per combination is measured with 100 repetitions on the held-out 25% test data. In addition, the standard deviation and the training accuracy are denoted. The idea is that, by using the ten best combinations, the results are not overly optimistic but represent a reasonable performance that could be reproduced with an average hyperparameter search. While it is theoretically possible that the found hyperparameters are biased towards the 75% training data, the evaluation on the independent test data should not yield over-optimistic results.

MoDT is used with three experts ($e = 3$) and a DT depth of two ($d = 2$). Derived from earlier experiments, this configuration seems to be an appropriate compromise between expressiveness and simplicity. Since DTs are well-known for their interpretability, and MoDT builds upon them, MoDT is compared to DTs with a maximum depth $d$ between two and four. Since $d = 4$ can already be problematic for interpretability (see Section V), a greater depth is not tested. As the Python library scikit-learn [17] is used for the DTs within the implementation of MoDT, the DTs for the comparison are also created with scikit-learn. Furthermore, MoDT is compared to RFs, which are valued for their good performance [1] and are also based on DTs. Although RFs are not interpretable, the comparison allows assessing MoDT from a pure performance perspective. Firstly, RFs are used with a configuration of $d = 2$ and $e = 3$ that is identical to the here used configuration of MoDT. Secondly, the default configuration of scikit-learn is used, which employs 100 DTs and imposes no restriction on the depth. The accuracy is evaluated using a 75/25% training/test split and 100 repetitions.

To allow a fair comparison, the same preprocessing steps that are included in MoDT are also applied to the training data of the alternatives. However, no extensive hyperparameter searches are conducted for the DTs and random forests. Yet, it is assumed that the default scikit-learn hyperparameters already yield reasonable results, especially in the deployed low complexity configurations. In addition, it would be problematic to compare the default DTs within the MoDT ensemble with optimized individual DTs. The complete above process is firstly performed for an MoDT configuration using the assumingly more interpretable but less powerful 2D gating function, and secondly for the full gating (FG) function.

*B. Results*

The average accuracies and the standard deviations using the MoDT classifier model with the 2D and FG function are erported in Table II. The best-performing interpretable methods on test accuracy (without consideration of the standard deviation) are highlighted with a grey background. RFs (as they are not interpretable) and the training accuracy of MoDT are not considered for this highlighting. The best overall methods on test accuracy are highlighted in bold.

In comparison to DTs with $d$ up to four, MoDT using a 2D gating function achieves a better test accuracy on 6 out of 8 datasets, hereby neglecting the standard deviations. In all cases, the test accuracy is close to the training accuracy (within 1–4 percentage points) except for the students dataset, where the test accuracy is significantly worse than the training accuracy. On this dataset, the DTs collectively outperform MoDT. In comparison with MoDT using the 2D gating function, the RFs' configuration with the same complexity as MoDT achieves lower scores on all datasets. Some of the differences are relatively large, e.g., MoDT performs 15 percentage points better on the steel dataset. MoDT achieves equal or better performance than the significantly more complex RF configuration with 100 DTs in 4 of the 8 datasets.

The full gating function exploits all features of $\mathbf{X}$ instead of just two. Is therefore can be more expressive but it is less interpretable. Therefore, it should be investigated, how this variant compares to the 2D gating variant. Using the full gating function, MoDT achieves a higher performance than the DTs in 6 out of 8 datasets. In contrast to the 2D variant, MoDT using a full gate now outperforms the DTs on the cars dataset. The opposite is observable for the contraceptive dataset, here, a DT with $d = 4$ now beats the performance of MoDT with a full gate. The full gate variant of MoDT achieves equal or better performance than the random forests configuration with 100 DTs in 5 of the 8 datasets.

Although the full gating function is in theory more powerful than its 2D counterpart, a direct comparison suggests that the full gating variant is not universally better. In 5 of the 8 datasets, both variants achieve roughly similar performance ($\pm 2$ percentage points). The full gating variant achieves a 10 percentage points better performance on the cars dataset. Conversely, the contraceptive and students datasets benefit from the 2D gate variant. Here, the test results are 5 and 4,

Table II

PERFORMANCE OF MODT WITH FG AND 2D GATE, DTs AND RFs. BEST INTERPRETABLE METHOD: GREY BACKGROUND, BEST OVERALL: BOLD.

| | MoDT 2D d=2 e=3 | | MoDT FG d=2 e=3 | | DT | | | RF | |
|---|---|---|---|---|---|---|---|---|---|
| dataset | training | test | training | test | d=2 | d=3 | d=4 | d=2 e=3 | d=* e=100 |
| abalone | .74±.00 | .71±.01 | .75±.01 | **.73±.01** | .67±.01 | .69±.01 | .70±.01 | .67±.01 | **.73±.01** |
| banknote | 1.00±.00 | .99±.00 | 1.00±.00 | **1.00±.00** | .91±.02 | .93±.01 | .95±.02 | .89±.04 | .99±.00 |
| breast | .96±.00 | .94±.01 | .97±.01 | **.95±.02** | .91±.02 | .92±.02 | .92±.02 | .91±.02 | .94±.02 |
| cars | .82±.01 | .78±.01 | .92±.01 | .88±.01 | .78±.02 | .79±.01 | .81±.02 | .71±.02 | **.96±.01** |
| contraceptive | .59±.01 | **.58±.01** | .57±.01 | .53±.02 | .48±.02 | .52±.03 | .55±.02 | .46±.04 | .52±.02 |
| iris | .99±.01 | .95±.02 | .99±.01 | **.96±.02** | .94±.03 | .94±.04 | .94±.03 | .92±.06 | .95±.03 |
| steel | .68±.01 | .68±.01 | .70±.02 | .67±.01 | .53±.02 | .53±.02 | .61±.02 | .53±.02 | **.78±.02** |
| students | .58±.02 | .45±.01 | .53±.01 | .41±.03 | .48±.04 | .51±.03 | .50±.04 | .44±.06 | .49±.03 |

respectively, percentage points higher. The training time of both variants is almost identical and ranges between $0.1\,\mathrm{s}$ for the iris dataset and $1.4\,\mathrm{s}$ for the abalone dataset.

## V. DISCUSSION

The comparison with interpretable DTs and RFs shows that MoDT is working as intended and can provide a benefit over individual DTs. On average, MoDT clearly outperforms single DTs with $d = 2$. MoDT continues to outperform DTs with $d = 4$, which are arguably already rather complex and can be hard to interpret. A DT of $d = 4$ can use more nodes (31) than MoDT with $d = 2$ and $e = 3$ (21). Therefore, it can be assumed that the gating function works as intended and cannot be replaced by a small number of additional DT splits.

On average, MoDT is clearly superior to a RF with the same complexity. As expected, the RF variant with a vastly more complex structure mostly outperforms MoDT. However, the resulting RFs are also not interpretable.Therefore, considering the way simpler structure of MoDT, it can be seen as a success that an interpretable method can achieve similar and sometimes even higher scores.

The comparison of the 2D gate and the full gate reveals that in all but one cases, namely the cars dataset, a 2D gating function is sufficient. This is an important result as the 2D gate is highly beneficial for interpretability. The variant can achieve similar and sometimes even better performance than the full gate. Possibly, the reduced complexity can prevent overfitting. However, the fixed number of iterations used in the experiment might be disadvantageous for the full gating variant. Since the variant is more complex, there are possibly more directions for optimization, and the training algorithm might need more iterations to explore them.

The most important hyperparameters of MoDT are the maximum depth $d$ of the DTs and the number of experts $e$. While it might seem difficult to set these values, a conceptional analysis of interpretability suggests a rather small range of suitable values. In one of the most cited papers in the field of psychology [18], the hypothesis known as *Miller's law* is introduced, which says that the average human can only keep approximately seven different chunks in their working memory. According to this, a binary DT with $d = 3$ (and thus up to 15 nodes) can already be difficult to understand. If such DTs are used within an ensemble, the comprehensibility becomes additionally difficult. DTs with $d = 1$ (decision stumps)

are naturally very comprehensible, however, the limited expressiveness diminishes their usefulness. Therefore, DTs of $d = 2$—thus, a maximum number of seven nodes—seem advisable. Regarding the number of experts, we assume that 2–5 experts are best suited for interpretability. Alternatively, it is also computationally feasible to simply test multiple configurations of MoDT with different numbers for $d$ and $e$.

If an MoDT model is used for prediction, potentially as a surrogate model of a more complex ML model like an artificial neural network, every decision can be retraced by a human. This allows experts to decide whether the prediction rules make sense. Additionally, in scenarios where potential misclassifications are dangerous, experts can rule out the possibility that the model will output certain predictions.

MoDT usually outperforms a DT with comparable complexity regarding prediction accuracy. Still, the model can remain completely comprehensible for humans if the 2D gating function is selected. The separation into smaller subregions reduces the cognitive load for a human. However, in edge cases, the visualization of the gating function can potentially lead to inaccuracies when it is used to retrace a decision of the model, e.g., when a data point is close to the border of two regions. Although this scenario is rare, the users of MoDT are implicitly enforced to incorporate this potential inaccuracy in their analysis.

It can be seen as a disadvantage for interpretability that the generated insights are not necessarily valid for the complete data, but only for distinct subregions. It might instead seem more desirable to have an analysis that generalizes to the complete dataset, e.g., by training just a single DT on the complete data. However, when the complexity of a DT must be limited such that it can be comprehensible by humans, the resulting DT might represent the real-world scenario poorly. In fact, a single short DT might be misleading as it can overly generalize complex relationships. Humans are prone to accept just one or two reasons as the sole explanations of possibly complex problems [19]. Therefore, a single short DT might lead humans to a rather incorrect analysis. Instead, with MoDT, a user is forced to reflect on multiple explanations for multiple subproblems. Possibly, this is a more realistic scenario for real-world problems.

MoDT has similarities to the commonly used IML method LIME [6]. Both methods can be used for explaining decisions

in local subregions. Yet, MoDT also enables global explanation as it does not need to fit a separate surrogate model for each observation. Further, MoDT's explanations are stable, i.e., identical inputs lead to identical explanations.

Generally, MoDT is not suited for all types of datasets and problems, respectively. If a dataset can be represented appropriately by a single DT, it does not make much sense to use MoDT instead of a DT. When MoDT is used on such a problem, it usually reduces itself to a single DT, regardless of the specified number of experts. Conversely, if an individual DT of arbitrary complexity seems entirely unsuited for a dataset, i.e., when a tree-based method cannot learn any meaningful pattern, MoDT is likely also not a good choice for the dataset. MoDT excels when a problem can be divided into multiple subproblems that incorporate different input-output relationships. Here, MoDT can outperform regular DTs as its gating function is more powerful and more flexible.

Although DTs are usually among the first examples of interpretable methods in IML literature, it is often not precisely defined how a DT can enable interpretability. While there are many implementations and extensions of DTs, most of them do not consider interpretability but purely focus on performance. We believe that there is potential for the improvement of DT algorithms and associated visualization methods. The algorithms should not treat interpretability as a byproduct but should be designed specifically with interpretability in mind. As MoDT builds upon DTs, MoDT would directly benefit from any improvements.

While there are numerous ways to improve MoDT from a technical perspective, we suggest that further work should initially focus on evaluating the method's benefit for interpretability in practice. This could be done with user surveys [20] or in-depth case studies. Ideally, the studies compare MoDT with regular DTs and a wide range of other IML methods [1]. For this, we believe that many existing studies wrongly reduce interpretability to single-dimensional properties like comprehensibility, trust, correctness or usefulness; but in practice, interpretability should be seen as the combination of multiple and sometimes competing objectives.

## VI. CONCLUSION

To the best of our knowledge, MoDT is the first ready-to-be-used DT ensemble method based on the MoE architecture that is specifically tailored for interpretability. The method employs a linear gating function that divides the input into disjoint subregions that are solved by distinct DTs. A novel approach is introduced that optionally limits the complexity of the gating function and thereby facilitates interpretability. Internally, the implementation uses a variation of the EM algorithm that employs a linear regression for finding model parameters instead of a more complex gradient-based optimization method which is typically found in MoE architectures. Enabling interpretability, MoDT provides a method to visualize the gating function which is adjusted to match the plots of the DTs. The experiments show that the implementation works and consistently outperforms individual DTs with similar complexity. Notably, the especially interpretable and simplistic 2D gating function can often compete with the more complex full gating function.
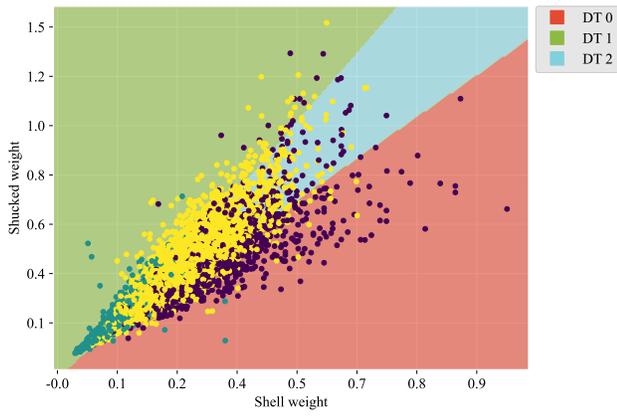
## REFERENCES

[1] C. Molnar, *Interpretable Machine Learning*, 2021. [Online]. Available: https://christophm.github.io/interpretable-ml-book

[2] W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller, *Explainable AI: interpreting, explaining and visualizing deep learning*. Springer Nature, 2019, vol. 11700.

[3] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, "Machine Learning Interpretability: A Survey on Methods and Metrics," *Electronics (Switzerland)*, vol. 8, 2019.

[4] T. Laugel, M. J. Lesot, C. Marsala, X. Renard, and M. Detyniecki, "The Dangers of Post-hoc Interpretability: Unjustified Counterfactual Explanations," in *28th International Joint Conference on Artificial Intelligence*, 2019, pp. 2801–2807.

[5] Z. C. Lipton, "The Mythos of Model Interpretability," *Queue*, vol. 16, no. 3, pp. 31–57, 2018.

[6] M. T. Ribeiro, S. Singh, and C. Guestrin, ""Why Should I Trust You?": Explaining the Predictions of Any Classifier," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016, p. 1135–1144.

[7] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[8] N. Burkart and M. Huber, "A Survey on the Explainability of Supervised Machine Learning," *Journal of Artificial Intelligence Research (JAIR)*, vol. 70, pp. 245–317, Jan. 2021.

[9] M. Vasic, A. Petrovic, K. Wang, M. Nikolic, R. Singh, and S. Khurshid, "MÖET: Interpretable and Verifiable Reinforcement Learning via Mixture of Expert Trees," *arXiv preprint*, 2019. [Online]. Available: https://arxiv.org/abs/1906.06717v3

[10] R. A. Jacobs, M. I. Jordan, S. E. Nowlan, and G. E. Hinton, "Adaptive Mixture of Experts," *Neural Computation*, vol. 3, pp. 79–87, 1991.

[11] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.

[12] Y. Yang and J. Ma, "A Single Loop EM Algorithm for the Mixture of Experts Architecture," *International Symposium on Neural Networks*, vol. 5552, pp. 959–968, 2009.

[13] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed. Springer, 2009.

[14] D. Dua and C. Graff, "UCI Machine Learning Repository," 2017. [Online]. Available: http://archive.ics.uci.edu/ml

[15] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.

[16] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A Next-generation Hyperparameter Optimization Framework," in *Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.

[17] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

[18] G. A. Miller, "The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information." *Psychological Review*, vol. 63, no. 2, pp. 81–97, 1956.

[19] T. Miller, "Explanation in Artificial Intelligence: Insights from the Social Sciences," *Artificial Intelligence*, vol. 267, pp. 1–38, 2019.

[20] R. Piltaver, M. Luštrek, M. Gams, and S. Martinčić-Ipšić, "What Makes Classification Trees Comprehensible?" *Expert Systems with Applications*, vol. 62, pp. 333–346, 2016.
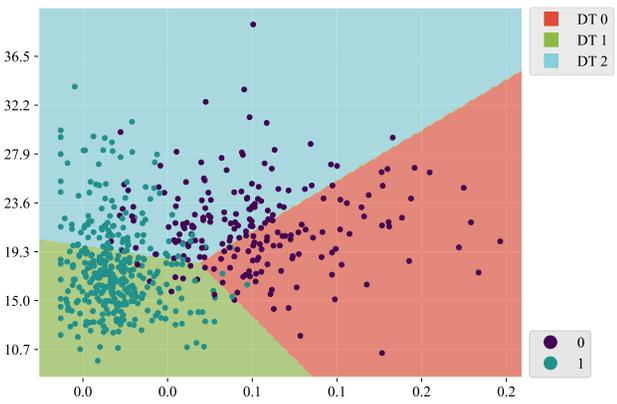
In Fig. 7 the 2D gating functions of MoDT with $d = 2$ and $e = 3$ for all datasets listed in Table I is plotted.
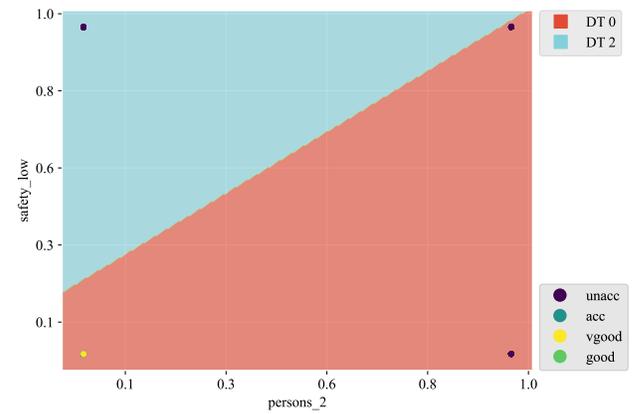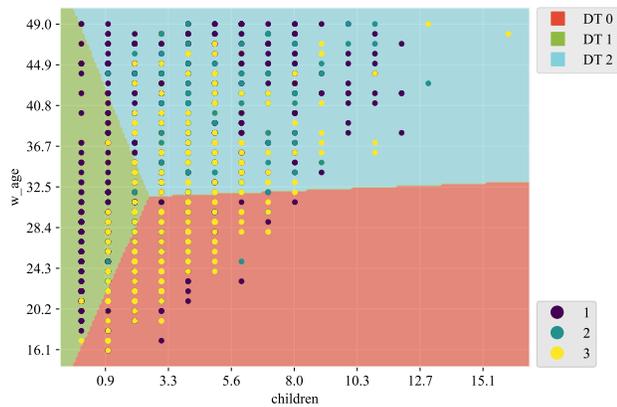


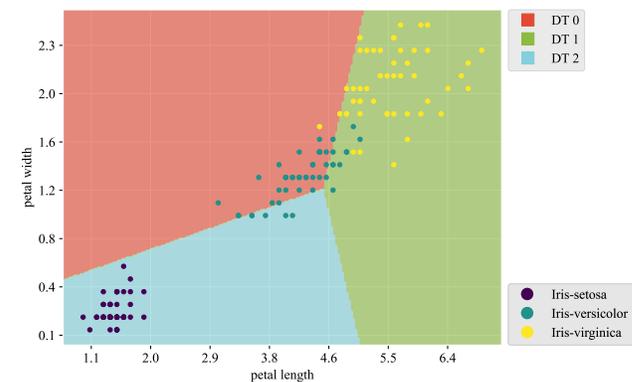(a) Abalone dataset.

(b) Banknote dataset.

(c) Breat dataset.

(d) Cars dataset. It is worth mentioning that the dataset only incorporates two values for the selected gating dimensions.
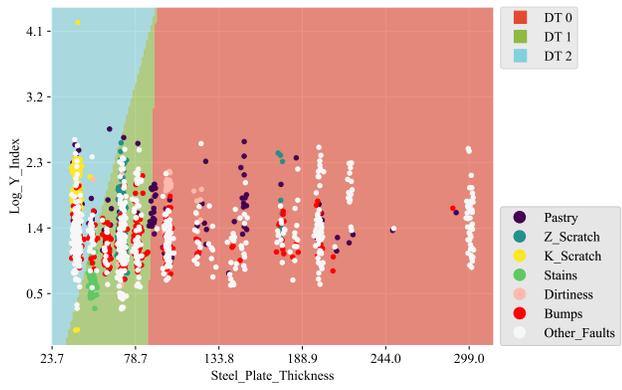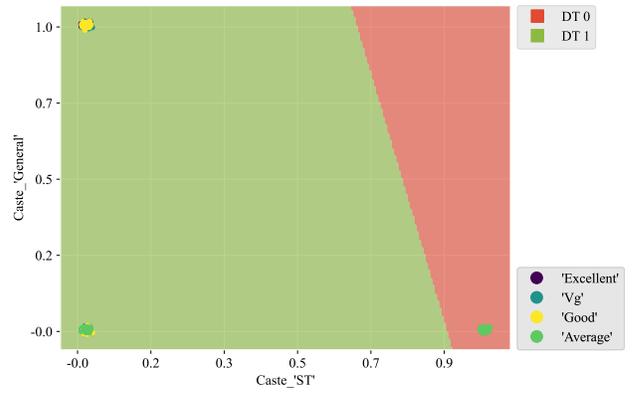
(e) Contraceptive dataset.

(f) Iris dataset.

Figure 7. 2D gating functions of MoDT.

(g) Steel dataset.

(h) Students dataset. It is worth mentioning that the dataset only incorporates two values for the selected gating dimensions.

Figure 7. 2D gating functions of MoDT.