# Interpretable Regional Descriptors:
# Hyperbox-Based Local Explanations

Susanne Dandl[1,2] , Giuseppe Casalicchio[1,2] , Bernd Bischl[1,2] , and Ludwig Bothmann[1,2]

[1] Department of Statistics, LMU Munich, Ludwigstr. 33, 80539 Munich, Germany
[2] Munich Center for Machine Learning (MCML), Munich, Germany
`Ludwig.Bothmann@stat.uni-muenchen.de`

**Abstract.** This work introduces interpretable regional descriptors, or IRDs, for local, model-agnostic interpretations. IRDs are hyperboxes that describe how an observation's feature values can be changed without affecting its prediction. They justify a prediction by providing a set of "even if" arguments (semi-factual explanations), and they indicate which features affect a prediction and whether pointwise biases or implausibilities exist. A concrete use case shows that this is valuable for both machine learning modelers and persons subject to a decision. We formalize the search for IRDs as an optimization problem and introduce a unifying framework for computing IRDs that covers desiderata, initialization techniques, and a post-processing method. We show how existing hyperbox methods can be adapted to fit into this unified framework. A benchmark study compares the methods based on several quality measures and identifies two strategies to improve IRDs.

**Keywords:** Interpretable machine learning, · Model-agnostic local interpretability · Semi-factual explanations · Hyperboxes.

## 1 Introduction

Supervised machine learning (ML) models are widely used due to their good predictive performance, but they are often difficult to interpret due to their complexity. Post-hoc interpretation methods from the field of interpretable machine learning (IML) can help to draw conclusions about the inner processes of these models. Two types of interpretation methods can be differentiated: local methods that explain individual predictions, and global methods that explain the expected behavior of the model in general. Doshi-Velez and Kim [4] define model interpretability as "the ability to explain or to present in understandable terms to a human". A topological form that satisfies this notion of interpretability is a hyperbox. In this work, we investigate hyperboxes as local interpretations that describe how the feature values of an observation can be changed without affecting its prediction. We call these boxes interpretable regional descriptors (IRDs). IRDs describe feature spaces by intervals for real-valued features and subsets of all possible classes for categorical features (see Table 1).

arXiv:2305.02780v1 [stat.ML] 4 May 2023

Table 1: Example based on the credit dataset [5,9] with 9 *features*. The second column shows the values of a *customer* with a moderate credit risk prediction. The *IRD* (generated by MaxBox & post-processing (Section 4)) shows how all features could be changed simultaneously so that the credit is still of moderate risk. The *1-dim IRD* shows how a single feature could be changed without changing the prediction (keeping the other features fixed). For features in the upper half, the IRD covers the full observed value *range* in the training data.

| Feature | Customer | IRD | 1-dim IRD | Range |
|---|---|---|---|---|
| sex | female | {female, male} | {female, male} | {female, male} |
| saving.accounts | little | {little, moderate rich} | {little, moderate, rich} | {little, moderate, rich} |
| purpose | car | {car, radio/TV, furniture, others} | {car, radio/TV, furniture, others} | {car, radio/TV, furniture, others} |
| age | 22 | [19, 22] | [19, 75] | [19, 75] |
| job | skilled | {skilled, highly skilled} | {unskilled, skilled, highly skilled} | {unskilled, skilled, highly skilled} |
| housing | rent | {rent} | {own, free, rent} | {own, free, rent} |
| checking.account | moderate | {little, moderate} | {little, moderate} | {little, moderate, rich} |
| credit.amount | 4000 | [4000, 5389] | [2127, 8424] | [276, 18424] |
| duration | 30 | [26, 33] | [6, 44] | [6, 72] |

## 1.1   Motivating Example for the Use of IRDs

Consider bank lending as a motivating example: a customer applies for a credit of €4000 at a bank to buy a new car. She is 22 years old, skilled, lives in a rented accommodation, has few savings and a moderate balance on her checking account. An ML model predicts whether the credit is of low, moderate or high risk. Due to a moderate risk prediction, the bank rejects the application. The IRD in Table 1 answers the question "to what extent the feature or multiple features can be changed such that the prediction is still in the moderate risk class". From an IRD, multiple insights into a prediction can be obtained.

First, IRDs offer a set of semi-factual explanations (SFEs) — also called a fortiori arguments — to justify a decision in the form of "even if" statements [19]. For these statements to be convincing, domain knowledge is required, e.g., that higher balances in the savings and checking account, and that higher skilled jobs decrease the risk for a bank. Given such knowledge, a multitude of semi-factual explanations can be derived from the IRD of Table 1 that (1) justify that a person is in the moderate risk class instead of the low risk class (e.g., "even if you had a moderate balance in the savings account and become highly skilled, your credit is still of moderate risk"), and that (2) justify that a person is in the moderate risk class instead of the high risk class ("even if you only have little balance in your checking account, your credit would still be of moderate risk"). The latter statement also reveals a "safety bound" if some of the features change in the undesirable direction (high risk class) in the future.

Second, the interval width or cardinality of a feature in an IRD relative to its entire feature space can indicate whether a feature affects a prediction locally (if Theorems 1 and 2 hold). For example, compared to the credit amount or balance status of her checking account, savings or purpose seem to have no local effect on the prediction in the bank lending example, since the regional descriptor covers the whole observed feature range in these two dimensions. These insights also reveal what can be options to change a given prediction.[3]

Third, IRDs are tools for model auditing. If the insights from a box (e.g., a semi-factual explanation) agree with domain knowledge, users have more trust in the model, while disagreement helps to reveal unintended pointwise biases or implausibilities of a model. For example, an IRD that does not cover male customers *might* indicate that the model classifies individuals differently based on gender.[4] An IRD that covers a credit amount of €300 and high balances in the checking account could be an indicator of an inaccurate model because such customers should pose only a low risk to the bank. In addition to credit risk, we show other practical applications of IRDs in Appendix A.

### 1.2   Contributions

Our contributions are: 1) We introduce IRDs as a new class of local interpretations to describe regions in the feature space that do not affect the prediction of an observation; 2) We formalize the search for IRDs as an optimization problem and develop desired properties of IRD methods; 3) We introduce a unifying framework for computing IRDs including initialization techniques and post-processing methods; 4) We show how existing hyperbox methods from data mining or IML can be adapted to fit into our unified framework; 5) We present a set of quality measures and compare our derived methods accordingly in a benchmark study; 6) We provide open-access repositories with an R package for the implemented approaches and the code for replicating the benchmark study.[5]

## 2   Methodology

Let $\hat{f} : \mathcal{X} \to \mathbb{R}$ be the prediction function of an ML model, where $\mathcal{X}$ denotes a $p$ dimensional feature space. For classification models, we consider a pre-defined class of interest for which $\hat{f}$ returns the predicted score or probability.

### 2.1   Formalizing the General Task for IRDs

Our goal is to find the largest hyperbox $B$ covering a point of interest $\mathbf{x}' \in \mathcal{X}$ where all data points in $B$ have a sufficiently close prediction to $\hat{f}(\mathbf{x}')$. The

---

[3] However, the concrete strategies can only reveal counterfactual explanations [27].

[4] Note that if all genders are part of the box, it does not mean the model is fair.

[5] Links will be shared upon acceptance. For review, we attached a zip file.

hyperbox $B$ should have $p$ dimensions $B = B_1 \times ... \times B_p$

$$\text{with } B_j = \begin{cases} \{c | c \in \mathcal{X}_j\} & \text{categorical } X_j \\ [l_j, u_j] \subseteq \mathcal{X}_j & \text{numeric } X_j \end{cases}, \tag{1}$$

consisting of intervals for numeric features and a subset of possible classes for categorical features. $\mathcal{X}_j$ reflects the value space of the $j$th feature $X_j$. In accordance with Lemhadri et al. [18], a prediction is sufficiently close if it falls into a *closeness region*, which is a user-defined prediction interval $Y' = [\hat{f}(\mathbf{x}') - \epsilon_L, \hat{f}(\mathbf{x}') + \epsilon_H]$ with $\epsilon_L, \epsilon_H \in \mathbb{R}_{\geq 0}$.[6] In the bank lending example, the closeness region should cover all model predictions that lead to the moderate risk class, e.g., a predicted probability of 30-60 % of defaulting, i.e., $Y' = [0.3, 0.6]$. To operationalize the above goal, we need three measures [22,24]:

1. $coverage(B) = \mathbb{P}(\mathbf{x} \in B | \mathbf{x} \in \mathcal{X})$, which measures how much a hyperbox covers the entire feature space. Since, in practice, not all $\mathbf{x} \in \mathcal{X}$ are observable, we use an empirical approximation given data $(\mathbf{x}_i)_{1 \leq i \leq n}$ with $\mathbf{x}_i \in \mathcal{X}$

$$\widehat{coverage}(B) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(\mathbf{x}_i \in B). \tag{2}$$

2. $precision(B) = \mathbb{P}(\hat{f}(\mathbf{x}) \in Y' | \mathbf{x} \in B)$, the fraction of points within a box $B$ whose predictions are inside $Y'$. Again, we use an empirical approximation

$$\widehat{precision}(B) = \frac{\sum_{i=1}^{n} \mathbb{I}(\mathbf{x}_i \in B \wedge f(\mathbf{x}_i) \in Y')}{\sum_{i=1}^{n} \mathbb{I}(\mathbf{x}_i \in B)}. \tag{3}$$

3. an indicator of whether $B$ covers $\mathbf{x}'$

$$locality(B) = \mathbb{I}(\mathbf{x}' \in B). \tag{4}$$

The following operationalizes the search for an IRD [22]:[7]

$$\begin{aligned} &\underset{B \subseteq \mathcal{X}}{arg\,max}(\widehat{coverage}(B)) \\ &\text{s.t. } \widehat{precision}(B) = 1 \text{ and } locality(B) = 1. \end{aligned} \tag{5}$$

**Definition 1.** *A box is maximal if and only if no box could be added under full precision, such that for all numeric $X_j$, it holds that $(\nexists x_j \in \mathcal{X}_j \wedge x_j < l_j : precision(B \cup [x_j, l_j]) = 1) \wedge (\nexists x_j \in \mathcal{X}_j \wedge x_j > u_j : precision(B \cup [u_j, x_j]) = 1)$, and for all categorical $X_j$, it holds that $(\nexists x_j \in \mathcal{X}_j \setminus B_j : precision(B \cup x_j) = 1)$.*

---

[6] For classification models, $Y' \subset [0, 1]$ must hold.

[7] For this, we extended the optimization task of Ribeiro et al. [22] to target IRDs by aiming for a precision of 1 and by including the locality constraint.

A box $B$ with maximum coverage satisfies this maximality property. We aim for a maximal $B$, since $B$ can then detect features that are not locally relevant for a prediction $\hat{f}(\mathbf{x}')$. We prove the following in Appendix B.

**Theorem 1.** *If $B$ is maximal, $B_j = [min(\mathcal{X}_j), max(\mathcal{X}_j)]$ holds for a feature $X_j$ that is not involved in the model $\hat{f}$.*

Similarly, we aim for homogeneous boxes $B$ such that $precision(B) = 1$. Then, $B$ can detect features that are locally relevant for $\hat{f}(\mathbf{x}')$. We prove the following in Appendix C.

**Theorem 2.** *If $precision(B) = 1$, $B_j \subset \mathcal{X}_j$ holds for a feature that is locally relevant for $\hat{f}(\mathbf{x}')$.*

## 2.2   Desiderata for IRDs

In Section 3, we discuss related methods to generate $B$. The suitability of these methods as IRD methods relies on whether they consider all objectives of Eq. (5) and whether they satisfy the following desired properties for IRDs.

*Interpretability*  In order for $B$ to be interpretable, we only consider methods that return a *single $p$-dimensional* hyperbox. The hyperrectangular structure of $B$ allows for a natural interpretation, which is not the case for hyperellipsoids or polytopes formed by halfspaces [18]. According to Eq. (5), $B$ needs to cover $\mathbf{x}'$, which is the case if the following holds: $\forall j \in \{1, ..., p\} : x'_j \in B_j$.

*Model-agnosticism*  The definition of $\hat{f}$ does not pose any restrictions on the ML model or the feature space. Therefore, methods should be model-agnostic such that they could explain both regression or classification models with various feature types (binary, nominal, ordinal or continuous).

*Sparsity constraints*  Eckstein et al. [6] proved that the optimization task for the maximum box problem is $\mathcal{NP}$-hard if the features defining the box are not fixed. This also applies to the search for IRDs, which only additionally requires $\mathbf{x}' \in B$. Since the search space for hyperboxes grows with the number of features, it is infeasible to consider all potential solutions. Furthermore, the fact that IRDs have as many dimensions as the dataset impedes their interpretability – the very goal of IRDs in the first place. To reduce the number of features, methods should be able to adhere to user-defined sparsity constraints such that for some features $X_j$, $B_j = x'_j$. Section 7 discusses other solutions.

## 3   Related Work

The optimization task of Eq. (5) can be understood mathematically as finding the preimage of prediction values $\in Y'$ in the neighborhood of $\mathbf{x}'$. Therefore,

Table 2: Overview of approaches that search for hyperboxes in feature spaces.

| | Objectives | | | Desiderata | | |
|---|---|---|---|---|---|---|
| | Coverage | Precision | Locality | Interpretable | Agnostic | Sparse |
| **Level set methods** | | | | | | |
| PBnB [28,29] | $\checkmark$ | $\checkmark$ | $\times$ | $\times$ | $\checkmark$ | $\times$ |
| **Data mining** | | | | | | |
| MaxBox [6] | $\checkmark$ | $\checkmark$ | $\times$ | $\checkmark$ | $\times$ | $\times$ |
| PRIM [10] | $\times$ | $\times$ | $\times$ | $\checkmark$ | $\times$ | $\times$ |
| **Post-hoc IML** | | | | | | |
| Anchors [22] | $\checkmark$ | $\checkmark$ | $\checkmark$ | $\checkmark$ | $\times$ | $\times$ |
| MAIRE [24] | $\checkmark$ | $\checkmark$ | $\checkmark$ | $\checkmark$ | $\times$ | $\times$ |
| LORE [12,13] | $\times$ | $\times$ | $\checkmark$ | $\checkmark$ | $\checkmark$ | $\times$ |
| **Interpretable classifier** | | | | | | |
| Column generation [2] | $\checkmark$ | $\checkmark$ | $\times$ | $\times$ | $\checkmark$ | $\times$ |

IRDs can be seen as a subset of a level set for function values $\in Y'$. Level set approximations often consist of points [7], and only a few approaches approximate these via hyperboxes [28,29]. These methods produce multiple boxes instead of a single one and do not require to contain a given $\mathbf{x}'$. Hence, they are not interpretable in our sense and, therefore, not useful to produce IRDs.

In data mining, [6] proposed a maximum box (MaxBox) approach for datasets with binary outcomes to find the largest homogeneous hyperbox w.r.t. the positive class. Friedman and Fisher [10] derived the Patient Rule Induction Method (PRIM) for seeking boxes in the feature space in which the outcome mean is high. Both approaches do not require $\mathbf{x}'$ to be in the box.

As described earlier, IRDs may also be seen as a method to summarize a multitude of SFEs. Most proposed methods for SFEs return only a single point as an explanation [3,14,19]. In contrast, the approach by Guidotti et al. [12,13] returns a set of SFEs using surrogate trees. Their approach reveals which feature values are most important for deriving a prediction by following the path to the point of interest. The reliability of surrogate trees depends on the assumption that the tree can adequately replicate the underlying model, which is often not the case. Furthermore, IRDs require homogeneous boxes, which is only possible with overfitting/deep-grown trees. Therefore, the tree structure is only suitable for deriving SFEs when the underlying model is tree-based [8,25].

An IML method that utilizes hyperboxes is the Anchors approach [22]. The returned hyperbox indicates how features must be fixed or anchored to prevent a model from changing the classification of a data point. Anchors were originally proposed to aim for hyperboxes that also partly cover observations of other classes; a precision of 0.95 is the default in its implementation [23]. Although the precision can be changed to 1, Anchors are nevertheless not suitable for the generation of IRDs due to their limited search space: Either the box boundary of a feature is set to the full feature range observed in the data, or to the value of $\mathbf{x}$. This bears the risk of "overly specific anchors" with low coverage [22]. To generate boxes with larger coverage, features can be binned beforehand.

However, no established discretization technique for Anchors exists so far and the optimization procedure underlying Anchors does not allow adaptions of the bins during optimization.

To overcome the discretization problem, Sharma et al. [24] proposed the model-agnostic interpretable rule extraction (MAIRE) procedure. MAIRE finds more optimal boundaries for continuous features via gradient-based optimization. It still does not allow a more precise choice for categorical features; either the box allows no changes to a feature or it covers all possible values of a feature.

Dash et al. [2] proposed a classifier based on a set of hyperboxes. The method focuses on an optimal combination of hyperboxes to derive an accurate model for inputs from the whole feature space using column generation. As such, the method does not focus on locality and is not interpretable in our sense.

Table 2 summarizes whether the addressed methods are suitable for generating IRDs. Overall, none of the methods satisfies all objectives of Eq. (5) and desiderata from Section 2.2. Specifically, none of them addresses sparsity constraints, and only a few are model-agnostic. In Section 4.4, we modify MaxBox, PRIM, and MAIRE such that they fulfill all of our requirements to transform them into useful IRD methods. All other methods cannot be modified to the required extent, since their underlying, irreplaceable optimization methods either target multiple boxes or different search spaces. The latter applies in particular to Anchors. However, the method serves as a baseline method for our benchmark study in Section 6.

## 4   Generating IRDs

We now present a unifying framework for generating IRDs, which consists of four steps: restriction, selection, initialization, and optimization. Optionally, a post-processing step can be conducted (Section 4.5).

### 4.1   Restriction of the Search Space

To restrict the initial search space for $B$, we propose a simple procedure to find the largest local box $\bar{\underline{B}}$ of $\mathbf{x}'$ such that $B \subset \bar{\underline{B}}$. For a continuous feature $X_j$, we vary its value $x_j'$ of $\mathbf{x}'$ on an equidistant grid. Upper and lower bounds of $\bar{B}_j$ are set to the minimal changes in $x_j'$, yielding a prediction outside $Y'$. This approach is similar to individual conditional expectation (ICE) values [11]. For a categorical feature $X_j$, $\bar{\underline{B}}_j$ comprises all classes of $\mathcal{X}_j$ that still lead to a prediction $\in Y'$ after adapting $x_j'$ of $\mathbf{x}'$. If a user sets the sparsity constraint that feature $X_j$ is immutable, $\bar{B}_j = x_j'$ must hold. We prove the following in Appendix D.

**Theorem 3.** *For any box $B$ that solves the optimization problem of Eq. (5) it holds that $B \subseteq \bar{B}$.*

## 4.2   Selection of the Underlying Dataset

All methods need a dataset $\bar{\mathbf{X}}$ consisting of $\mathbf{x} \in \mathcal{X}$ as an input. This dataset is used for evaluating (competing) boxes w.r.t. the empirical versions of coverage and precision (Eq. (2) and Eq. (3)). For some methods, the dataset also offers a set of potential box boundaries to be evaluated. A suitable dataset is the training data. Since only instances $\in \bar{\mathrm{B}}$ are relevant (Theorem 3), we remove all instances $\notin \bar{\mathrm{B}}$ from $\bar{\mathbf{X}}$. Consequently, $x_j = x'_j \, \forall \, \mathbf{x} \in \bar{\mathbf{X}}$ holds for all immutable features $X_j$. More features and sparsity constraints increase the risk that $\bar{\mathbf{X}}$ is only sparsely populated around $\mathbf{x}'$. Since we aim for IRDs that are faithful to the model and not to the data-generating process (DGP), new data can be generated by uniformly sampling from the admissible feature ranges of $\bar{\mathrm{B}}$. In Section 6, we inspect how double-in-size sampled data within $\bar{\mathrm{B}}$ [8] affects the quality of IRDs and IRD methods compared to using training data.

## 4.3   Initialization of a Box

All methods require an initial box $B$ as an input, which is either set to the largest local box $\bar{\mathrm{B}}$ covering all $\bar{\mathbf{X}}$ or the smallest box possible which only contains $\mathbf{x}'$. We define methods that start with the largest local box as top-down IRD methods, and methods that start with the smallest box possible as bottom-up methods.

## 4.4   Optimization of Box Boundaries

The last step comprises the optimization of the box boundaries. Top-down methods iteratively shrink the box boundaries of the largest local box to improve the box's precision (upholding that $\mathbf{x}' \in B$), while bottom-up methods iteratively enlarge the box boundaries of the smallest box to improve the box's coverage (upholding the precision at 1). In this section, we describe the MaxBox, MAIRE, and PRIM approaches and our extensions such that the methods optimize Eq. (5) and fulfill the desiderata of Section 2.2. Pseudocodes and illustrations of the inner workings of the extended approaches are given in Appendix E. All methods receive as input a dataset $\bar{\mathbf{X}}$ and an initial box $B$.

*MaxBox – Top-down Method* MaxBox was originally proposed for binary classification problems – with a positive and negative class. The method starts with the largest box covering all data. A branch and bound (BnB) algorithm [17] inspects the options to shrink the box to optimize its precision w.r.t. the positive class. The branching rule creates new boxes by bracketing out a sample $\mathbf{x}$ of the negative class, such that the box is shrunk to be either below or above the values of $\mathbf{x}$ in at least one feature dimension (categorical features are one-hot encoded). Estimates of the upper bound for the coverage of a box determine which imprecise box is branched next, which sample is used for branching, and which boxes are discarded because their upper bound does not exceed the coverage of

---

[8] Double-in-size refers to the size of the training data, not of $\bar{\mathbf{X}}$.

the current largest homogeneous box. If no boxes to shrink are left, the largest homogeneous box is returned as an IRD.

*Extensions* By labeling observations with predictions $\in Y'$ as positive, the approach becomes model-agnostic. Since the original algorithm does not consider whether corresponding boxes still include $\mathbf{x}'$, we adapted the approach to discard boxes that do not contain $\mathbf{x}'$ to guarantee locality.

*PRIM – Top-down Method* The method originally aims for boxes with a high average outcome. The procedure starts with a box that includes all points. In the peeling phase, PRIM iteratively identifies a set of eligible subboxes (defined by the $\alpha$- and $(1\text{-}\alpha)$-quantile for numeric features and each present category for categorical features) and peels off the subbox that results in the highest average outcome after exclusion. This step is repeated until the number of points included in the box drops below a fraction of the total number of points. In the pasting phase, the box is iteratively enlarged by adding the subbox that increases the outcome mean the most. These subboxes consist of at least $\alpha$ observations with the nearest lower or higher values in one dimension (numeric $X_j$) or with a new category (categorical $X_j$).

*Extensions* We adapted the approach to target Eq. (5): in each peeling iteration, the subbox is excluded such that the resulting box has the highest precision (coverage acts as a tiebreaker), and in each pasting iteration, the largest homogeneous subbox is added. If the precision and coverage are not sufficient to select a best box for peeling or pasting, a subbox is randomly selected from the best ones. Peeling stops as soon as the resulting box is homogeneous, while pasting stops as soon as there exists no homogeneous box to add. Furthermore, only subboxes that do not cover $\mathbf{x}'$ are peeled. According to the authors' recommendation, we use $\alpha = 0.05$ for the benchmark study (Section 6).

*MAIRE – Bottom-up Method* The method starts with a box covering $\mathbf{x}'$. In each iteration, the box boundaries are adapted via ADAM [15] by optimizing a differentiable approximation of the coverage measure. If the precision falls below a certain threshold or $\mathbf{x}'$ is not part of the box, the method additionally optimizes a differentiable version of Eq. (3) and Eq. (4), respectively. MAIRE stops after a specified number of iterations. In the end, the method returns the largest homogeneous box over the iterations.

*Extensions* The method requires 0-1-scaled features. To overcome the one-vs-all issue for categorical features (Section 3), we one-hot-encode categorical features. We implemented a convergence criterion for a fair comparison with the other (convergent) approaches: we let MAIRE enlarge the box boundaries until the precision falls below 1, then MAIRE is only allowed to run for another 100 iterations. The implementation for the experiments in Section 6 is based on the authors' implementation [24] with the discussed modifications. The hyperparameters were set according to the authors' recommendations. We only set the precision threshold to 1, rather than 0.95.

### 4.5   Post-processing

All methods described in the previous section determine box boundaries based on a finite number of data points in $\bar{\mathbf{X}}$. The limited access carries the risk that some regions of the feature space are not represented in $\bar{\mathbf{X}}$ and that the boundaries of a generated $B$ are suboptimal: There could be areas in $B$ that have predictions $\notin Y'$, or there could be adjacent areas outside of $B$ that also have predictions $\in Y'$. To improve the box boundaries of a given box $B$, we developed the following post-processing method using newly sampled data. The procedure consists of peeling and pasting as PRIM.

First, the precision of $B$ is measured based on newly sampled data. If $\exists \mathbf{x} \in B$ with $\hat{f}(\mathbf{x}) \notin Y'$, subboxes with the lowest precision in proportion to their size (according to newly sampled data within this subbox) are iteratively peeled. If all subboxes to peel are homogeneous, peeling stops. In the subsequent pasting step, the largest subboxes that proved to be homogeneous (according to newly sampled data within this subbox) are added. If the best box cannot be clearly determined (because several boxes have the same precision and coverage), a subbox is randomly chosen. The method has three hyperparameters: the number of samples used for evaluation, the relative box size (in relation to the size of $\mathcal{X}_j$) for peeling or pasting boxes for continuous features, and a threshold for the minimum box size. The latter acts as a stopping criterion for pasting. If no homogeneous subbox can be added, the relative box size to add for continuous features is halved as long as the relative box size is not lower than the threshold. The pseudocode of our method displays Appendix F.

Section 6 investigates whether our post-processing method improves IRDs. For the experiments, we set the number of samples to evaluate boxes to 100, the relative box size to 0.1, and the threshold for the minimum box size to 0.05.

## 5   Quality Measures

We now present a set of quality measures for *generated IRDs* and *IRD methods*. These measures apply to a single instance $\mathbf{x}'$ to be explained, where $B$ is the returned IRD of $\mathbf{x}'$ of an IRD method $G$. The assessment requires evaluation data $\mathbf{E}$ consisting of $\mathbf{x} \in \mathcal{X}$; for the benchmark study in Section 6, we use training data and new data uniformly sampled from $\bar{\mathrm{B}}$. Training data helps to assess whether the methods use the training data appropriately during IRD generation (e.g., precision should be 1), while a proliferated number of newly generated data $\in \bar{\mathrm{B}}$ leads to a more precise evaluation w.r.t. the model, not the DGP.

*Locality* The IRD should cover $\mathbf{x}'$. This property is fulfilled if $locality(B) = \mathbb{I}(\mathbf{x}' \in B)$ equals 1.

*Coverage* Given two IRDs with equal precision, we prefer the one with higher coverage (Eq. (2)). To evaluate the coverage, we use samples $\mathbf{x} \in \mathbf{E}$ from the connected convex level set $\mathcal{L}$ covering $\mathbf{x}'$.

**Definition 2.** *A datapoint* $\mathbf{x}$ *with* $\hat{f}(\mathbf{x}) \in Y'$ *is part of* $\mathcal{L}$ *of* $\mathbf{x}'$ *iff there exists a path between* $\mathbf{x}$ *and* $\mathbf{x}'$ *for which all intermediate points have a prediction* $\in Y'$.

Paths are identified via the identification algorithm of Kuratomi et al. [16], details are given in Appendix G.

*Precision* Given two IRDs with equal coverage, the IRD with higher precision is preferred (Eq. (3)).

*Maximality* A box should be maximal according to Definition 1 based on $\mathbf{x} \in \mathbf{E}$ instead of $\mathbf{x} \in \mathcal{X}$.

*No. of Calls* Lower number of calls to $\hat{f}$ of an IRD method are preferred.

*Robustness* If we rerun method $G$ on the same $\mathbf{x}'$ and $\hat{f}$ $R$ times using the same $\bar{\mathbf{X}}$, the produced IRDs $B_1, ..., B_R$ should overlap with the originally produced $B$, such that $robustness(G) = \min\limits_{k \in \{1,...,R\}} \frac{\sum_{\mathbf{x} \in \mathbf{E}} \mathbb{I}(\mathbf{x} \in B \cap B_k)}{\sum_{\mathbf{x} \in \mathbf{E}} \mathbb{I}(\mathbf{x} \in B \cup B_k)}$ has a high value.

## 6 Performance Evaluation

In a benchmark study, we address the following research questions (RQs):

1. Based on the stated quality measures of Section 5, how do the different methods of Section 4.4 perform against each other and the baseline method when training data are used as $\bar{\mathbf{X}}$ (without post-processing)?
2. What effect do double-in-size sampled data originating from $\bar{B}$ have on the quality of the IRDs and methods compared to using training data?
3. What effect does the post-processing (Section 4.5) have on the quality of the IRD methods?

As a baseline method, we use the Anchors approach [22] with a precision of 1 and 20-quantile-based bins for numeric features (see Section 3 for details).

### 6.1 Setup

To answer the RQs, we utilize six datasets available on the OpenML platform [26], either with a binary, multi-class or continuous target variable. Table 3 summarizes the datasets' dimensions as well as the target and feature types. Before training a model, five randomly sampled datapoints were excluded from the datasets to be $\mathbf{x}'$. On each of the datasets, four models are trained: a hyperbox model, a logistic regression/multinomial/linear model (depending on the outcome), a neural network with one hidden layer, and a random forest model. The number of trees for the random forest model and the neurons on the hidden layer are tuned (details are given in Appendix H). The hyperbox model is derived from a classification and regression tree (CART) model for each $\mathbf{x}'$ individually. For a

Table 3: Overview of benchmark datasets. ID: OpenML id; Type: target type; Obs: number of rows; Cont/Cat: number of continuous/categorical features.

| Name | ID | Type | Obs | Cont | Cat |
|---|---|---|---|---|---|
| diabetes | 37 | binary | 768 | 8 | 0 |
| tic_tac_toe | 50 | binary | 958 | 0 | 9 |
| cmc | 23 | three-class | 1473 | 2 | 7 |
| vehicle | 54 | four-class | 846 | 18 | 0 |
| no2 | 886 | regression | 500 | 7 | 0 |
| plasma_retinol | 511 | regression | 315 | 10 | 3 |

given $\mathbf{x}'$, the post-processed model predicts 1 if a point falls in the same terminal node as $\mathbf{x}'$ and 0 otherwise.[9]

For classification models, the prediction function returns the probability of the class with the highest probability for $\mathbf{x}$. For binary targets, we set $Y' = [0.5, 1]$. For regression and multi-class targets, $Y'$ is set to $[\hat{f}(\mathbf{x}) - \delta, \hat{f}(\mathbf{x}) + \delta]$ with $\delta$ as the standard deviation of predictions $\hat{f}$ of the training data. For multi-class, the interval is additionally capped between 0 and 1. For each dataset, model, and $\mathbf{x}'$, we generate IRDs with MaxBox, PRIM, and MAIRE, as well as Anchors – our baseline method. The hyperparameters of the methods were set according to Section 4. The methods were either run on training or on uniformly sampled data from $\bar{\mathrm{B}}$ (RQ 2), and either without or with post-processing (RQ 3). For the robustness evaluation, we repeated the experiments $R = 5$ times.

The methods and their generated IRDs were evaluated based on the performance measures of Section 6 – either evaluated on the training data or 1000 new instances sampled uniformly from $\bar{\mathrm{B}}$. We also compared the methods statistically by conducting Wilcoxon rank-sum tests for the hypothesis that the distribution of the coverage and precision values do not differ between two (IRD) methods (RQ 1), for a method using training vs. sampled data (RQ 2), and for a method without vs. with post-processing (RQ 3). The experiments were conducted on a computer with a 2.60 GHz Intel(R) Xeon(R) processor, and 32 CPUs. Overall, generating the boxes took 63 hours spread over 20 CPUs. The five repetitions for the robustness evaluation required another 316 hours.

### 6.2   Results

Figure 1 compares the coverage and precision values of the methods visually. Table 4 shows the frequency of fulfilling maximality and the number of calls to $\hat{f}$ of the methods. The separate results for each dataset and model, the statistical analysis, and the results of robustness are shown in Appendix I. We omitted the results for the locality measure because all returned IRDs covered $\mathbf{x}'$.

*RQ 1 - comparison of methods* Without post-processing and training data as $\bar{\mathbf{X}}$ (first row, Figure 1), MaxBox had the highest precision as evaluated on

---

[9] The true hyperbox of the CART model might be larger than the terminal node-induced hyperbox (see Figure 6 in the Appendix).

Table 4:  Comparison of methods w.r.t. maximality and no. of calls to $\hat{f}$ averaged over all datasets, models and $\mathbf{x}'$. Each method was run or evaluated on training data or uniformly sampled data from $\bar{\mathrm{B}}$, and without (0) or with (1) post-processing. Higher maximality and lower no. of calls are better.

| | Traindata | | | | | | Sampled | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mathrm{Max}_{train}$ | | $\mathrm{Max}_{samp}$ | | No. calls to $\hat{f}$ | | $\mathrm{Max}_{train}$ | | $\mathrm{Max}_{samp}$ | | No. calls to $\hat{f}$ | |
| | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| MaxBox | **0.60** | **0.42** | 0.06 | **0.41** | **184** | 55769 | 0.23 | **0.45** | 0.24 | **0.43** | **1621** | **37627** |
| PRIM | 0.42 | 0.37 | **0.18** | 0.39 | **184** | 46070 | 0.20 | 0.42 | **0.25** | 0.39 | **1621** | 42958 |
| MAIRE | 0.18 | 0.41 | 0.04 | **0.41** | **184** | 68126 | 0.06 | 0.41 | 0.11 | 0.35 | **1621** | 92976 |
| Anchors | 0.27 | **0.42** | 0.16 | 0.40 | 26402 | 94448 | **0.31** | 0.42 | 0.18 | 0.36 | 77818 | 129276 |

training and newly sampled data, followed by MAIRE. The IRDs of PRIM had on average the largest coverage, but they also covered sampled data with predictions outside $Y'$. Due to the randomized choice of a subbox in the case of ties, PRIM is not robust according to our robustness metric. None of the methods outperformed the other methods w.r.t. maximality. Overall, all methods outperformed the baseline method Anchors according to coverage and precision and calls to $\hat{f}$. The latter is because competing boxes are evaluated on column-wise permutations of the observed data. All other methods only called $\hat{f}$ $|\bar{\mathbf{X}}|$ times.

*RQ 2 - training vs. sampled data* On average, double-in-size sampled data originating from $\bar{\mathrm{B}}$ led to slightly higher coverage, precision and maximality rates w.r.t. newly sampled data but not w.r.t. the training data. Due to the increase in the size of $\bar{\mathbf{X}}$, more calls to $\hat{f}$ were necessary.[10]

*RQ 3 - without vs. with post-processing* Post-processing increased the coverage and precision of IRDs for all methods. The difference in the quality of IRDs between the methods and between the underlying data scheme (training data vs. sampled data) diminished. Quality enhancement comes at the cost of efficiency and robustness; on average, post-processing resulted in 57,000 additional calls to $\hat{f}$ and the sampling of new data decreased the robustness. MAIRE required on average the most post-processing iterations, followed by Anchors.

## 7  Conclusion, Limitations and Outlook

*Conclusion* We introduced IRDs that describe regions in the feature space that do not affect the prediction of an instance in the form of hyperboxes. These hyperboxes provide a set of semi-factual explanations to justify a prediction, and indicate which features affect a prediction and whether there might be pointwise biases or implausibilities. We formalized the search for IRDs, and introduced desiderata, a unifying framework and quality measures for IRD methods. We

---

[10] The size decuples instead of doubles compared to the training data, because not all training data are $\in \bar{\mathrm{B}}$ and, thus, not in $\bar{\mathbf{X}}$.
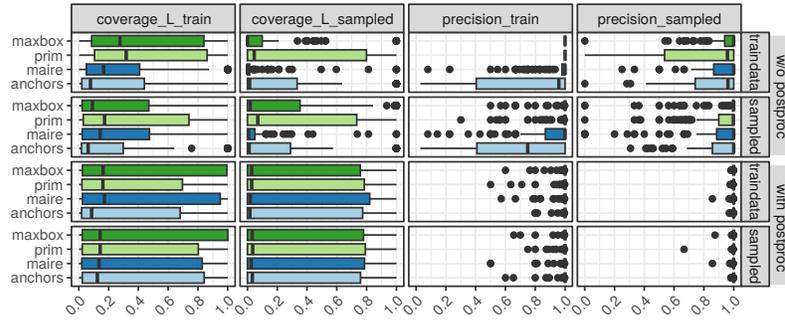
Fig. 1: Comparison of IRD methods w.r.t. coverage and precision as evaluated on the training data or newly sampled data within $\bar{\mathrm{B}}$. Addendum L means that for the coverage evaluation only training or sampled points within $\mathcal{L}$ are considered. Each point in the boxplot reflects the performance of a generated IRD of one experimental setting (dataset, model and $\mathbf{x}'$). Each method was either run or evaluated on training data (traindata) or uniformly sampled data from $\bar{\mathrm{B}}$ (sampled), and the methods were run either without or with post-processing (postproc). Higher values for precision and coverage are better.

discussed three existing hyperbox methods in detail and adapted them to search for IRDs. The lack of a method "ruling it all" in the benchmark study emphasizes the need for a unifying framework comprising multiple methods. The study also revealed that access to a larger, uniformly sampled dataset or using our proposed post-processing method can further enhance the quality of IRDs.

*Limitations* Our work offers potential for further research, e.g., on the sensitivity of the methods' hyperparameters, on the influence of sampling sizes, or on the methods' robustness w.r.t. slight changes in $\mathbf{x}'$ or the underlying data. While we only considered low-dimensional datasets in the benchmark study, for high-dimensional datasets we proposed two strategies to restrict the search space: either by letting users decide which features can be changed and which cannot (Section 2.2), or by deriving the largest local box $B \subset \bar{\mathrm{B}}$ based on ICE curves (Section 4.1). Further research can explore: (1) the use of other IML methods, such as feature importance methods, to select features for which changes are investigated (all other features are set to their admissible value range); (2) the consideration of feature correlations or causal relations to generate IRDs, which not only naturally restricts the search space but also makes the IRD faithful to the DGP. Considering feature correlations is also important for the application of IRDs beyond tabular data. While all presented methods are model-agnostic, we leave concrete investigations on image and text data to future research.

*Outlook* We believe that our work can also be a starting point for investigations on the application of IRDs in other fields, e.g., for hyperparameter (HP) tuning: if a promising HP set for an ML model was identified by a tuning method, IRDs can

reveal its sensitivity and whether there are other equally good but more efficient HP settings. IRDs might also identify high-fidelity regions for interpretable local surrogate models, like LIME [21]. LIME approximates predictions of a black-box model $\hat{f}(\mathbf{x})$ around an observation $\mathbf{x}'$ using a (regularized) linear model $\hat{g}(\mathbf{x})$. Here, it might be useful to understand in which region $B$ the linear model approximates the black-box model (high-fidelity region); $\hat{g}$ only provides valuable insights in the region $B$ around $\mathbf{x}'$ where $\forall \mathbf{x} \in B : \hat{h}(\mathbf{x}) := |\hat{f}(\mathbf{x}) - \hat{g}(\mathbf{x})| \leq \epsilon$ for a user-defined $\epsilon > 0$. With $\hat{h}$ as the prediction model and $Y' = [0, \epsilon]$, IRD methods might identify such high-fidelity regions $B$ in an interpretable manner.

## Acknowledgements

## References

1. Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: SMOTE: Synthetic minority over-sampling technique. Journal of Artificial Intelligence Research **16**, 321–357 (2002). https://doi.org/10.1613/jair.953
2. Dash, S., Günlük, O., Wei, D.: Boolean decision rules via column generation. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. p. 4660–4670. NIPS'18, Curran Associates Inc., Red Hook, NY, USA (2018)
3. Dhurandhar, A., Chen, P.Y., Luss, R., Tu, C.C., Ting, P., Shanmugam, K., Das, P.: Explanations based on the missing: Towards contrastive explanations with pertinent negatives. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. p. 590–601. NIPS'18, Curran Associates Inc., Red Hook, NY, USA (2018)
4. Doshi-Velez, F., Kim, B.: Towards a rigorous science of interpretable machine learning. arXiv 1702.08608 v2, arXiv.org E-Print Archive (2017). https://doi.org/10.48550/arXiv.1702.08608
5. Dua, D., Graff, C.: Uci machine learning repository (2017), https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data)
6. Eckstein, J., Hammer, P.L., Liu, Y., Nediak, M., Simeone, B.: The maximum box problem and its application to data analysis. Computational Optimization and Applications **23**(3), 285–298 (2002). https://doi.org/10.1023/a:1020546910706
7. Emmerich, M.T.M., Deutz, A.H., Kruisselbrink, J.W.: On quality indicators for black-box level set approximation. In: Tantar, E., Tantar, A.A., Bouvry, P., Del Moral, P., Legrand, P., Coello Coello, C.A., Schütze, O. (eds.) EVOLVE-A Bridge between Probability, Set Oriented Numerics and Evolutionary Computation, pp. 157–185. Springer Berlin Heidelberg, Berlin, Heidelberg (2013). https://doi.org/10.1007/978-3-642-32726-1_4
8. Fernandez, G., Aledo, J.A., Gamez, J.A., Puerta, J.M.: Factual and counterfactual explanations in fuzzy classification trees. IEEE Transactions on Fuzzy Systems **30**(12), 5484–5495 (2022). https://doi.org/10.1109/tfuzz.2022.3179582

9. Ferreira, L.: German credit risk (2018), https://www.kaggle.com/datasets/kabure/german-credit-data-with-risk, last accessed 23.01.2023

10. Friedman, J.H., Fisher, N.I.: Bump hunting in high-dimensional data. Statistics and Computing **9**(2), 123–143 (1999). https://doi.org/10.1023/A:1008894516817

11. Goldstein, A., Kapelner, A., Bleich, J., Pitkin, E.: Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. Journal of Computational and Graphical Statistics **24**(1), 44–65 (2015). https://doi.org/10.1080/10618600.2014.907095

12. Guidotti, R., Monreale, A., Giannotti, F., Pedreschi, D., Ruggieri, S., Turini, F.: Factual and counterfactual explanations for black box decision making. IEEE Intelligent Systems **34**(6), 14–23 (2019). https://doi.org/10.1109/MIS.2019.2957223

13. Guidotti, R., Monreale, A., Ruggieri, S., Naretto, F., Turini, F., Pedreschi, D., Giannotti, F.: Stable and actionable explanations of black-box models through factual and counterfactual rules. Data Mining and Knowledge Discovery (2022). https://doi.org/10.1007/s10618-022-00878-5

14. Kenny, E.M., Keane, M.T.: On generating plausible counterfactual and semi-factual explanations for deep learning. Proceedings of the AAAI Conference on Artificial Intelligence **35**(13), 11575–11585 (2021). https://doi.org/10.1609/aaai.v35i13.17377

15. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv 1412.6980 v9, arXiv.org E-Print Archive (2017). https://doi.org/10.48550/arXiv.1412.6980

16. Kuratomi, A., Miliou, I., Lee, Z., Lindgren, T., Papapetrou, P.: JUICE: JUstIfied Counterfactual Explanations. In: Pascal, P., Ienco, D. (eds.) Discovery Science. pp. 493–508. Lecture Notes in Computer Science, Springer Nature Switzerland, Cham (2022). https://doi.org/10.1007/978-3-031-18840-4_35

17. Land, A.H., Doig, A.G.: An automatic method of solving discrete programming problems. Econometrica **28**(3), 497–520 (1960). https://doi.org/0.2307/1910129

18. Lemhadri, I., Li, H.H., Hastie, T.: RbX: Region-based explanations of prediction models. arXiv 2210.08721, arXiv.org E-Print Archive (2022). https://doi.org/10.48550/arXiv.2210.08721

19. Nugent, C., Doyle, D., Cunningham, P.: Gaining insight through case-based explanation. Journal of Intelligent Information Systems **32**(3), 267–295 (2009). https://doi.org/10.1007/s10844-008-0069-0

20. Pfisterer, F., Poon, J., Lang, M.: mlr3keras: mlr3 keras extension. Github repository. URL https://github.com/mlr-org/mlr3keras (2022),
Commit: bad8434b7898b51b2143fc680594057c00dc7080

21. Ribeiro, M.T., Singh, S., Guestrin, C.: "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. pp. 1135–1144 (2016)

22. Ribeiro, M.T., Singh, S., Guestrin, C.: Anchors: High-precision model-agnostic explanations. Proceedings of the AAAI Conference on Artificial Intelligence **32**(1) (2018). https://doi.org/10.1609/aaai.v32i1.11491

23. Ribeiro, M.T., Singh, S., Guestrin, C.: Anchor. Github repository. URL https://github.com/marcotcr/anchor (2022),
Commit: b1f5e6ca37428613723597e85c38558e8cd21c2e

24. Sharma, R., Reddy, N., Kamakshi, V., Krishnan, N.C., Jain, S.: MAIRE - a model-agnostic interpretable rule extraction procedure for explaining classifiers. In: Holzinger, A., Kieseberg, P., Tjoa, A.M., Weippl, E. (eds.) Machine Learning and Knowledge Extraction, vol. 12844, pp. 329–349. Springer International

Publishing, Cham (2021). https://doi.org/10.1007/978-3-030-84060-0_21, Series Title: Lecture Notes in Computer Science

25. Stepin, I., Alonso, J.M., Catala, A., Pereira-Fariña, M.: Generation and evaluation of factual and counterfactual explanations for decision trees and fuzzy rule-based classifiers. In: 2020 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE). pp. 1–8. IEEE, Glasgow, United Kingdom (2020). https://doi.org/10.1109/FUZZ48607.2020.9177629

26. Vanschoren, J., van Rijn, J.N., Bischl, B., Torgo, L.: OpenML: Networked science in machine learning **15**(2), 49–60 (2014). https://doi.org/10.1145/2641190.2641198

27. Wachter, S., Mittelstadt, B., Russell, C.: Counterfactual explanations without opening the black box: Automated decisions and the GDPR. Harvard Journal of Law & Technology **31**(2), 841–887 (2018)

28. Zabinsky, Z.B., Huang, H.: A partition-based optimization approach for level set approximation: Probabilistic branch and bound. In: Smith, A.E. (ed.) Women in Industrial and Systems Engineering: Key Advances and Perspectives on Emerging Topics, pp. 113–155. Springer International Publishing, Cham (2020). https://doi.org/10.1007/978-3-030-11866-2_6

29. Zabinsky, Z.B., Wang, W., Prasetio, Y., Ghate, A., Yen, J.W.: Adaptive probabilistic branch and bound for level set approximation. In: Proceedings of the 2011 Winter Simulation Conference (WSC). pp. 4146–4157. IEEE, Phoenix, AZ, USA (2011). https://doi.org/10.1109/WSC.2011.6148103

## A    Application Examples

In addition to the credit application in Section 1, we show in the following a medical and jurisdictional application.

*Medical* Consider an ML model that predicts if a person will develop diabetes in the future. (For simplicity, we assume this model accurately approximates real world relationships.) In the following, we discuss two cases:

(1) A person that is predicted to develop diabetes wants to know why this is the case and what can be options to prevent this. There are different potential actions to take: more sport, less red meat, homeopathic medicine, etc. The IRD can tell which action is not promising, e.g., sports when all realistic amounts of sport are inside the box. However, changing the diet might be an option, because changing the diet by just eating meat one day a week is not part of the box (concrete strategies for prevention can reveal counterfactual explanations).

(2) A person that is predicted not to develop diabetes wants to know how flexible their life-style is without changing the prediction. It may be okay for a person to gain weight without having a higher risk of developing diabetes, as long as they do not change their diet towards including more red meat.

*Jurisdiction* Consider an ML model that predicts if a person will commit a crime in the next 2 years. A person that gets a high score wants to know why. IRDs that do not contain all groups of protected attributes, such as gender, can indicate unfair discrimination against these groups. Hence, IRDs can initiate further investigations on fairness and biases of an ML model.

## B   Proof of Theorem 1

*Proof.* Given a feature $X_j$ that is not involved in the prediction model $\hat{f}$ such that $\forall \tilde{\mathbf{x}} \in \mathcal{X} \wedge \forall x_j \in \mathcal{X}_j$:

$$\hat{f}(\tilde{x}_1, ..., \tilde{x}_{j-1}, \tilde{x}_j, \tilde{x}_{j+1}, ..., \tilde{x}_p) = \hat{f}(\tilde{x}_1, ..., \tilde{x}_{j-1}, x_j, \tilde{x}_{j+1}, ..., \tilde{x}_p), \qquad (6)$$

and given a box $B$ for $\mathbf{x}'$ that is maximal according to Definition 1. We assume now that Theorem 1 does not hold such that $B_j = [l_j, u_j] \subset \mathcal{X}_j$. However, since Eq. (6) holds, either $(\exists\, x_j \in \mathcal{X}_j \wedge x_j < l_j : precision(B \cup [x_j, l_j]) = 1)$, or $(\exists\, x_j \in \mathcal{X}_j \wedge x_j > u_j : precision(B \cup [u_j, x_j]) = 1)$ for numeric $X_j$ or $(\exists\, x_j \in \mathcal{X}_j \setminus B_j : precision(B \cup x_j) = 1)$ for categorical $X_j$ holds which contradicts the maximality assumption of $B$.

## C   Proof of Theorem 2

*Proof.* Given a box $B$ with $precision(B) = 1$ and $\mathbf{x}' \in B$, and given a feature $X_j$ that is relevant for $\hat{f}(x')$ such that $\exists x_j \in \mathcal{X}_j \setminus B_j : \hat{f}(x'_1, ..., x'_{j-1}, x_j, x'_{j+1}, ..., x'_p) \notin Y'$. We assume now that Theorem 2 does not hold, such that $B_j = \mathcal{X}_j$. This contradicts the statement that $precision(B) = 1$ because $x_j$ that leads to a prediction $\notin Y'$ for $\mathbf{x}'$ is also covered by the box.

## D   Proof of Theorem 3

*Proof.* Without loss of generality, we assume that we only have numeric features. Assume we computed $\bar{\mathrm{B}} = \bigcup_{j=1}^{p} [l_j, u_j]$ such that $\forall j \in \{1, ...p\}$ :

$$\hat{f}(\underbrace{x'_1, .., x'_{j-1}, l_j, x'_{j+1}, ..., x'_p}_{:=\mathbf{x}'_l}) \notin Y' \wedge \hat{f}(\underbrace{x'_1, .., x'_{j-1}, u_j, x'_{j+1}, ..., x'_p}_{:=\mathbf{x}'_u}) \notin Y'.$$

We assume that $B \subset \bar{\mathrm{B}}$ is not true for now such that there is a homogeneous $B$ with $min(B_j) < l_j$ or $max(B_j) > u_j$ and $\mathbf{x}' \in B$. However, then either $\mathbf{x}'_l$ or $\mathbf{x}'_u$ would also be part of $B$ but for both $\hat{f}(\mathbf{x}'_u) \notin Y'$ or $\hat{f}(\mathbf{x}'_l) \notin Y'$ holds, which contradicts that $B$ is homogeneous.

# E   Pseudocode and Illustrations of IRD Methods

## E.1   Pseudocode

---

**Algorithm 1** Adapted MaxBox approach [6]

---

**Input:** Targeted instance $\mathbf{x}'$, desired range $Y'$, prediction model $\hat{f} : \mathcal{X} \to \mathbb{R}$, input dataset $\bar{\mathbf{X}}$, initial box $B$

Initialize candidates $= [\,]$, upper_bound_coverage_best $=$ -Inf, current_best $= [\,]$

**if** $\exists \mathbf{x} \in \bar{\mathbf{X}} \wedge \mathbf{x} \in B : \hat{f} \notin Y'$ **then**

  candidates $=$ candidates.$append(B)$

  **while** $length$(candidates) $> 0$ **do**

    $B^{best} = choose\_best$(candidates)

        $\triangleright$ if upper_bound_coverage_best $< 0$, $B^{best}$ corresponds to the box with the most no. of shrinking steps done before (with the upper bound of the coverage as a tiebreaker), else, $B^{best}$ corresponds to the box that maximizes $\left( \frac{|\{\mathbf{x} \in B | \hat{f}(\mathbf{x}) \in Y'\}|}{|\{\mathbf{x} \in B | \hat{f}(\mathbf{x}) \notin Y'\}|} \right)$.

    candidates $=$ candidates.$remove(B^{best})$

    children $= create\_new\_candidates(B^{best})$   $\triangleright$ in Figure 2, C and D are new candidates created from the initial box

    **for** $B \in$ children **do**

      **if** $\forall \mathbf{x} \in B : \hat{f}(\mathbf{x}) \in Y'$ **then**

        coverage $= upper\_bound\_coverage(B)$

        **if** coverage $>$ upper_bound_coverage_best **then**

          current_best $= B$

          upper_bound_coverage_best $=$ coverage

        **end if**

      **else**

        **if** $upper\_bound\_coverage(B) >$ upper_bound_coverage_best **then**

          candidates $=$ candidates.$append(B)$

        **end if**

      **end if**

    **end for**

  **end while**

**else**

  current_best $=$ B

**end if**

**return** current_best

---

---

**Algorithm 2** Adapted PRIM approach [10]

---

**Input:** Targeted instance $\mathbf{x}'$, desired range $Y'$, prediction model $\hat{f} : \mathcal{X} \to \mathbb{R}$, input dataset $\bar{\mathbf{X}}$, initial box $B$

**while** $\exists \mathbf{x} \in \bar{\mathbf{X}} \wedge \mathbf{x} \in B : \hat{f} \notin Y'$ **do**

   **for** $j \in \{1, ..., p\}$ **do**

      $C_j = [\,]$                           ▷ create candidates for peeling

      **if** $X_j$ numeric **then**

         $C_j = C_j.append(B_j^-, B_j^+)$ where $B_j^- = [l_j, min(X_{j(\alpha)}, x_j')]$ and

$B_j^+ = [max(X_{j(1-\alpha)}, x_j'), u_j]$ with $x_{j(\alpha)}$ and $x_{j(1-\alpha)}$ as the $\alpha$- and $(1-\alpha)$-quantiles of $X_j$ in the current box $B$

      **else if** $X_j$ categorical **then**

         $C_j = \{s \in B_j \mid s \neq x_j'\}$

      **end if**

   **end for**

   $b^{best} = \underset{b \in C_j,\, j \in \{1,...,p\}}{\arg\max}\; precision(B \setminus b)$

   $B = B \setminus b^{best}$

**end while**

homogeneous = TRUE

**while** homogeneous **do**

   **for** $j \in \{1, ..., p\}$ **do**

      $C_j = [\,]$                           ▷ create candidates for pasting

      **if** $X_j$ numeric **then**

         inbox = $\{\mathbf{x} \in \bar{\mathbf{X}} \mid x_k \in B_k\}$, for $k \in \{1, ..., j-1, j+1, ...p\}$

         number_added = $|\{\mathbf{x} \in \bar{\mathbf{X}} \mid \mathbf{x} \in B\}| \cdot \alpha$

         $C_j = C_j.append(B_j^-, B_j^+)$ with $B_j^- = [x_j^l, l_j]$ and $B_j^+ = [u_j, x_j^u]$ with

$x_j^l$ as the $j$th feature value of the (number_added)th observation $\mathbf{x} \in$ inbox with a value $x_j$ lower than $l_j$ and

$x_j^u$ as the $j$th feature value of the (number_added)th observation $\mathbf{x} \in$ inbox with a value $x_j$ higher than $u_j$

      **else if** $X_j$ categorical **then**

         $C_j = \{s \in X_j \mid s \notin B_j\}$

      **end if**

      $C_j = \{b \in C_j \mid precision(B \cup b) = 1\}$

   **end for**

   **if** $\exists j \in \{1, ..., p\} : |C_j| > 0$ **then**

      $b^{best} = \underset{b \in C_j,\, j \in \{1,...,p\}}{\arg\max}\; coverage(B \setminus b)$

      $B = B \cup b$

   **else**

      homogeneous = FALSE

   **end if**

**end while**

**return** B

---

---

**Algorithm 3** Adapted MAIRE approach [24]

---

**Input:** Targeted instance $\mathbf{x}'$, desired range $Y'$, prediction model $\hat{f} : \mathcal{X} \to \mathbb{R}$, input dataset $\bar{\mathbf{X}}$, initial box $B$, precision threshold $\tau$ (default 1), maximum number of iterations max_iterations (default 100)

Scale all feature values of $\mathbf{x} \in \bar{\mathbf{X}}$ and $\mathbf{x}'$ to 0-1 range

best_coverage $= 0$

converged $=$ FALSE

best_candidate $= B$

$i = 0$

**while** $i \leq max\_iterations$ **do**

  $B = optimize\_with\_adam(B)$

      $\triangleright$ optimizes differentiable versions of coverage, precision and locality

  **if** $precision(B) \geq \tau \wedge coverage(B) \geq$ best_coverage **then**

    best_candidate $= B$

  **else if** $precision(B) < \tau$ **then**

    converged $=$ TRUE

  **end if**

  **if** converged $=$ TRUE **then**

    i $=$ i $+ 1$

  **end if**

**end while**

**return** best_candidate

---

### E.2 Illustrations



Fig. 2: Illustration of the adapted MaxBox algorithm. The algorithm starts with $\bar{\mathrm{B}}$ (dashed box). In the box are two data points with predictions $\notin Y'$ (called negative samples) and the box needs to be further optimized. First, a negative sample is chosen - either the one in A or B. Therefore, the number of samples with predictions $\in Y'$ after excluding the points in one feature dimension are inspected. The resulting boxes of both negative samples cover a maximum of seven samples. We chose the one of A (B is also fine). Its resulting boxes are the new subproblems/candidates (C and D). Both boxes in C and D only include samples with predictions $\in Y'$, but the box in C is chosen as an optimum because it includes more samples with predictions $\in Y'$. D is discarded because it has a lower number. Since C and D cannot be further split because no negative samples are within both boxes, the returned box by MaxBox is the box in C.
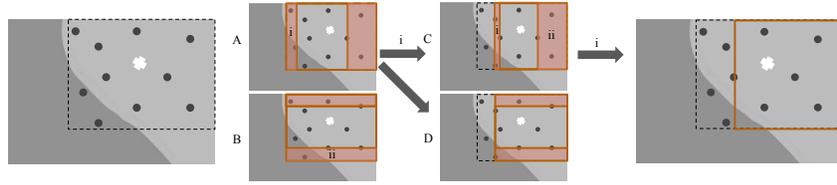
Fig. 3: Illustration of the adapted PRIM algorithm. The algorithm starts with $\bar{\mathrm{B}}$. In the first iteration, there exist four potential subboxes (two in each feature dimension (A vs. B)) that could be removed. The subbox i is chosen because it has the highest precision but compared to ii it has a smaller size. In the next step (C & D), again four subboxes can be potentially removed. Again, we choose i for the same reason as before. After its removal, the resulting box is at the same time the final box because in the pasting step only one subbox could be added – i again. All other dimensions are maximal.



Fig. 4: Illustration of the adapted MAIRE algorithm. The algorithm starts with the smallest box possible. The box boundaries are then iteratively enlarged (A-D). The box boundaries are only updated if the precision of the new box = 1.

## F   Pseudocode of post-processing Approach

---

**Algorithm 4** Post-processing algorithm - peeling (inspired by [10])

---

**Input:** Targeted instance $\mathbf{x}'$, desired range $Y'$, prediction model $\hat{f} : \mathcal{X} \to \mathbb{R}$, initial box $B$, number of samples for evaluation $M$ (default 100), relative subbox size of continuous features $\alpha$ (default 0.1)

**for** $j \in \{1, ..., p\}$ **do**
  **if** $X_j$ numeric **then**
    $s_j = (max(\mathcal{X}_j) - min(\mathcal{X}_j)) \cdot \alpha$         ▷ derive subbox sizes for numeric features based on $\mathcal{X}$
    **if** $X_j$ integer **then**
      $s_j = round(s_j)$
    **end if**
  **end if**
**end for**
$\bar{\mathbf{X}} = sample\_uniformly(B, n = M \cdot 5)$         ▷ sample new data to check if $B$ homogeneous
**if** $\exists \mathbf{x} \in \bar{\mathbf{X}} \land \mathbf{x} \in B : \hat{f} \notin Y'$ **then**
  not\_homogeneous = TRUE         ▷ start peeling
  **while** not\_homogeneous **do**
    **for** $j \in \{1, ..., p\}$ **do**
      $C_j = []$         ▷ create candidates for peeling
      **if** $X_j$ numeric **then**
        $C_j = C_j.append(B_j^-, B_j^+)$
where $B_j^- = [l_j, min(l_j + s_j, x'_j)]$ and $B_j^+ = [max(u_j - s_j, x'_j), u_j]$
      **else if** $X_j$ categorical **then**
        $C_j = \{s \in B_j \mid s \neq x'_j\}$
      **end if**
      $C_j = \{b \in C_j \mid precision(B_j^b) < 1\}$ with $B_j^b = (B_1 \times ... \times B_{j-1} \times b \times B_{j+1} \times ... \times B_p)$
    **end for**
    **if** $\exists j \in \{1, ..., p\} : |C_j| > 0$ **then**
    $b^{best} = \underset{b \in C_j, \, j \in \{1,...,p\}}{\arg\max} \; precision\_to\_boxsize(B_j^b)$   ▷ evaluate on $M$ new instances sampled within $B_j^b$
    $B^{best} = (B_1 \times ... \times B_{j-1} \times b^{best} \times B_{j+1} \times ... \times B_p)$ ▷ choose the one with lowest precision relative to size
      $B = B^{best}$
    **else**
      not\_homogeneous = FALSE
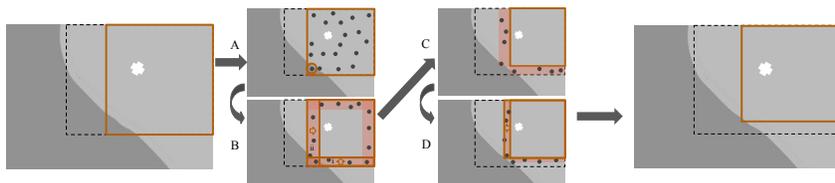    **end if**
  **end while**
**end if**
**return** B, $\mathbf{s} = \{s_j \mid X_j$ numeric$\}$

---

---

**Algorithm 5** Post-processing algorithm - pasting (inspired by [10])

---

**Input:** Targeted instance $\mathbf{x}'$, desired range $Y'$, prediction model $\hat{f} : \mathcal{X} \to \mathbb{R}$, initial box $B$ (potentially peeled), number of samples for evaluation $M$ (default 100), relative subbox size of continuous features $\alpha$ (default 0.1), lower threshold for relative subbox size $\alpha_0$ (default 0.05), subbox sizes of numeric features $\mathbf{s}$

homogeneous = TRUE                    $\triangleright$ start pasting

stepsize = 1

**while** homogeneous **do**

  **for** $j \in \{1, ..., p\}$ **do**

    $C_j = []$                    $\triangleright$ create candidates/subboxes for pasting

    **if** $X_j$ numeric **then**

      $C_j = C_j.append(B_j^-, B_j^+)$

where $B_j^- = [l_j - \text{stepsize} \cdot s_j, l_j]$ and $B_j^+ = [u_j, u_j + \text{stepsize} \cdot s_j]$

    **else if** $X_j$ categorical **then**

      $C_j = \{s \in X_j \mid s \notin B_j\}$

    **end if**

    $C_j = \{b \in C_j \mid precision(B_j^b) = 1\}$ with $B_j^b = (B_1 \times ... \times B_{j-1} \times b \times B_{j+1} \times ... \times B_p)$

  **end for**

  **if** $\exists j \in \{1, ..., p\} : |C_j| > 0$ **then**

    $b^{best} = \underset{b \in C_j, \, j \in \{1,...,p\}}{\arg\max} size(B_j^b)$    $\triangleright$ evaluate on $M$ new instances sampled within $B_j^b$

    $B = B \cup b$    $\triangleright$ choose largest one with precision 1

  **else**

    **if** stepsize $\geq \alpha_0$ **then**

      stepsize = stepsize/2                    $\triangleright$ if no box with precision 1 exists, consider reducing the subbox sizes

    **else**

      homogeneous = FALSE

    **end if**

  **end if**

**end while**

**return** B

---

Fig. 5: Illustration of the post-processing algorithm. The algorithm starts with the box generated by another method (solid brown box, which is a subbox of the dashed box $\bar{\underline{B}}$). First, new points are sampled and it is assessed whether the box is homogeneous (A). If not, the subboxes with the lowest precision compared to their size are peeled iteratively (B). The precision is assessed based on newly sampled points within the subboxes. First subbox i is peeled then subbox ii (both contain a sample with a prediction $\notin Y'$). If no subbox with precision $< 1$ exists, it is assessed whether the box could be further enlarged (C). If all considered subboxes have precisions $< 1$, the subbox sizes are halved (D) as long as the relative subbox size does not fall below a threshold.



## G   Level Set Identification

The algorithm starts at $\mathbf{x}'$ and tries to find a connection $\in Y'$ between the nominal, then the ordinal, and then the continuous features of $\mathbf{x}$ and $\mathbf{x}'$. If a path is found, $\mathbf{x}$ is part of $\mathcal{L}$. For categorical features, all permutations of feature orders are inspected.[11] For continuous features, the shortest linear path for a given number of equidistant steps is checked. Kuratomi et al. [16] used DBSCAN, for which the choice of the maximum distance threshold is ambiguous. The identification algorithm has a complexity of $O(c! \cdot c + o! \cdot \sum_{j=1}^{o} k_j + q)$ with $c$ and $o$ as the number of nominal and ordinal features, respectively, $k_j$ as the number of possible values of an ordinal feature $X_j$ and $q$ as the number of inspected steps for continuous features.

The level set could be further enriched by attempting to find connections between the unconnected and connected points. For the comparison of IRD methods, however, a convex level set is sufficient, since the hyperbox itself is convex.

## H   Tuning of ML models

For hyperparameter tuning, we used random search (with 15 evaluations), and 5-fold cross-validation (CV) with the misclassification error (classification) or mean squared error (regression) as a performance measure. Table 5 shows the tuning search space of each model. The rather limited tuning setup should be

---

[11] If the number of permutations exceeds 100 permutations, 100 feature orders are randomly chosen.

sufficient for our task of explaining a prediction model – a less accurate model is not a hindrance. Unbalanced datasets such as *tic_tac_toe*, *diabetes* and *cmc* were balanced with the SMOTE algorithm [1]. For SMOTE, numeric features were standardized and categorical ones were one-hot encoded. The optimizer for the neural network was ADAM [15] with 500 epochs. For all other hyperparameters, the default values of the mlr3keras R package were used [20] (apart from the no. of layer units, see Table 5). Table 6 shows the accuracies of each model using nested resampling with 5-fold CV in the inner and outer loop).

Table 5: Tuning search space of each model. Hyperparameter values of *num.trees* were log-transformed.

| Model | Hyperparameter | Range |
|---|---|---|
| random forest | num.trees | [1, 1000] |
| logistic regression | - | - |
| linear model | - | - |
| multi-nomial model | - | - |
| hyperbox/rpart | - | - |
| neural net | layer_units | [1, 20] |

Table 6: Classification error or mean squared error (regression) of each model on each dataset. The performances were computed using nested resampling with 5-fold CV in the inner and outer loop. We did not measure the performance of the (terminal node) hyperbox model because the model differs for each $\mathbf{x}'$.

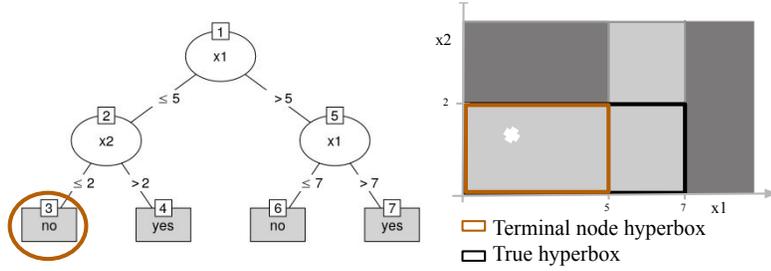| | Random forest | Linear model | Neural net | Hyperbox |
|---|---|---|---|---|
| diabetes | 0.233 | 0.224 | 0.229 | - |
| tic_tac_toe | 0.036 | 0.019 | 0.094 | - |
| cmc | 0.466 | 0.495 | 0.389 | - |
| vehicle | 0.256 | 0.201 | 0.254 | - |
| no2 | 33502.856 | 37678.319 | 77866.331 | - |
| plasma_retinol | 45391.218 | 59224.452 | 297481.249 | - |

Fig. 6: True hyperbox vs. terminal node hyperbox for a CART tree. The white cross corresponds to $\mathbf{x}'$.

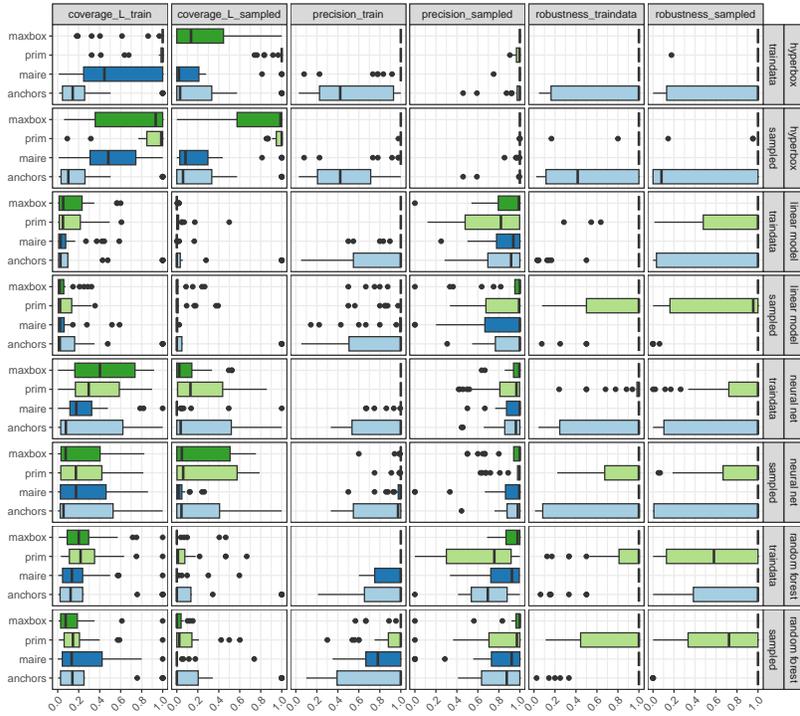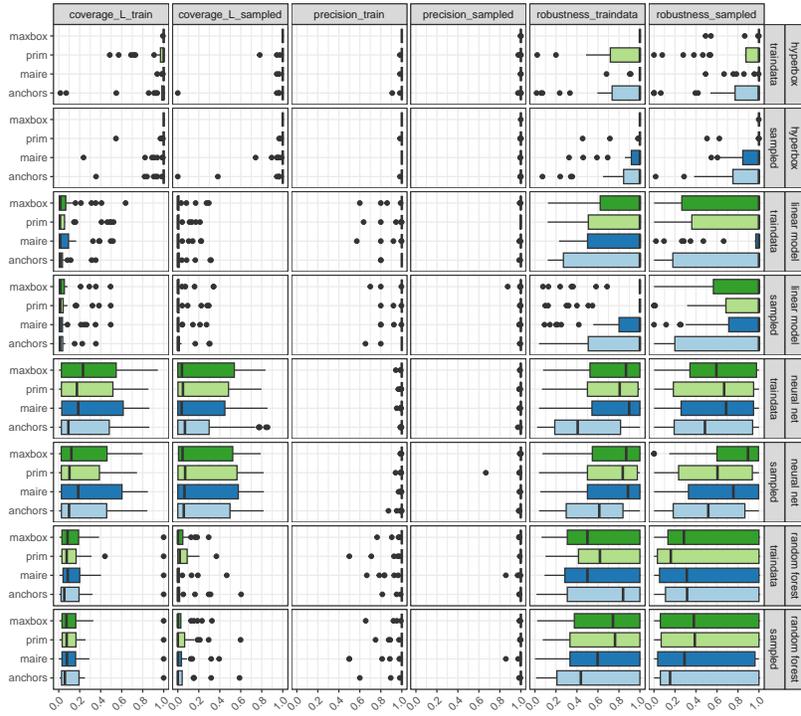# I   Benchmark - Additional Results



Fig. 7: Comparison of MaxBox, PRIM, Anchors, and MAIRE w.r.t. coverage and precision for each model separately. Each method was either run or evaluated on training data (traindata) or uniformly sampled data from $\bar{\mathrm{B}}$ (sampled) *without* post-processing. Higher values for precision and coverage are better.

Fig. 8: Comparison of MaxBox, PRIM, Anchors, and MAIRE w.r.t. coverage and precision for each model separately. Each method was either run or evaluated on training data (traindata) or uniformly sampled data from $\bar{B}$ (sampled) *with* post-processing. Higher values for precision and coverage are better.
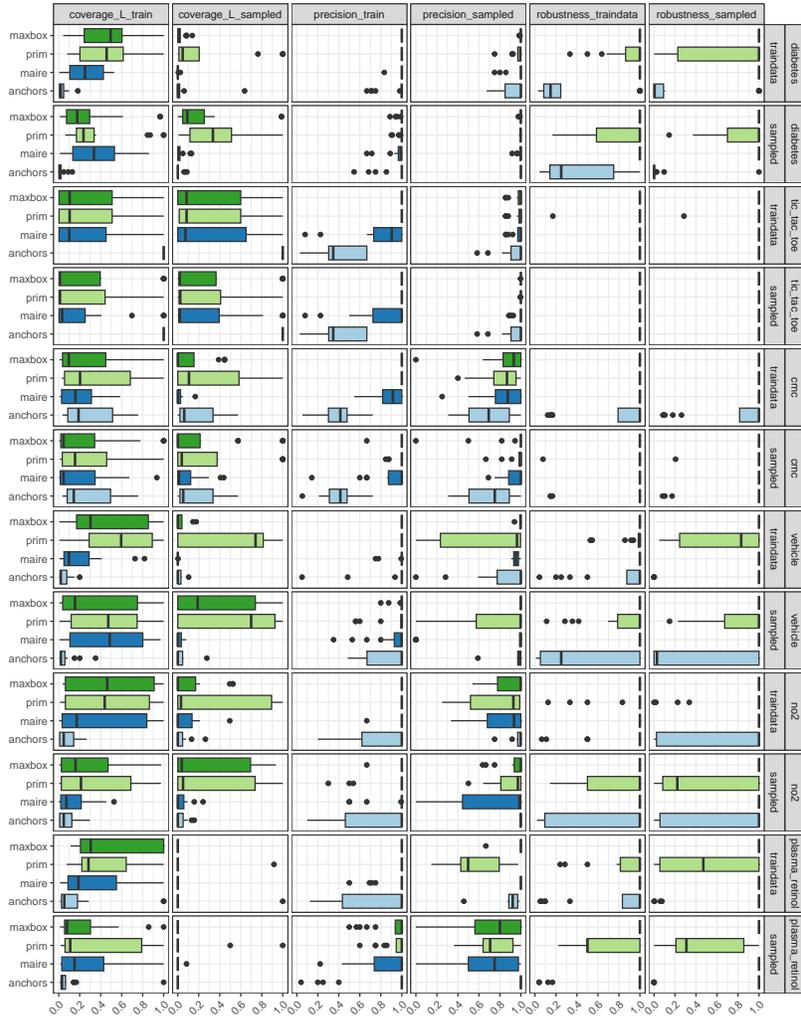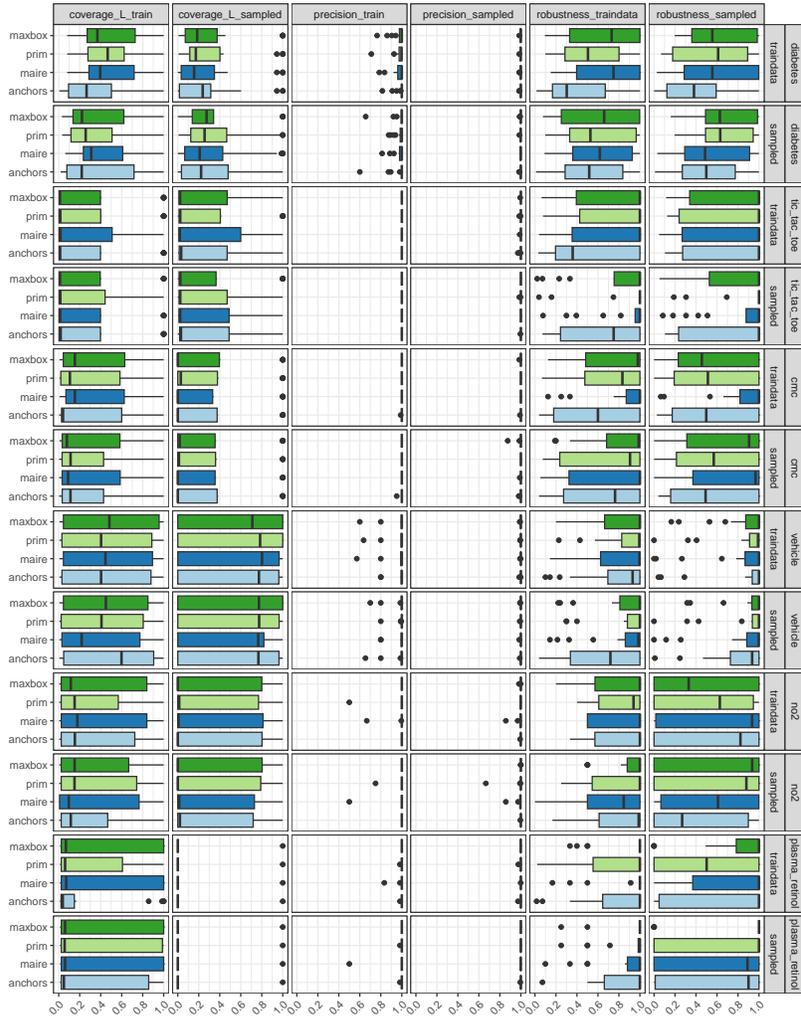
Fig. 9: Comparison of MaxBox, PRIM, Anchors, and MAIRE w.r.t. coverage and precision for each dataset separately. Each method was either run or evaluated on training data (traindata) or uniformly sampled data from $\bar{\mathrm{B}}$ (sampled) *without* post-processing. Higher values for precision and coverage are better.

Fig. 10: Comparison of MaxBox, PRIM, Anchors, and MAIRE w.r.t. coverage and precision for each dataset separately. Each method was either run or evaluated on training data (traindata) or uniformly sampled data from $\bar{\mathrm{B}}$ (sampled) *with* post-processing. Higher values for precision and coverage are better.

Table 7: Statistical analysis of RQ 1. Pairwise comparison of MaxBox, PRIM, Anchors, and MAIRE w.r.t. coverage and precision. Each value corresponds to the p-value obtained for the Wilcoxon rank-sum test with $H_0$ that the performances do not differ. Cells printed in bold font correspond to p-values that are lower than $\alpha = 0.05/36$ (Bonferroni-adjustment) and indicate that one method outperforms the other. Only methods run on the training data without post-processing were compared.

| measure | MaxBox = PRIM | MaxBox = Anchors | MaxBox = MAIRE | PRIM = Anchors | PRIM = MAIRE | Anchors = MAIRE |
|---|---|---|---|---|---|---|
| coverage_train | 0.761 | 0.618 | **0** | **0** | 0.579 | 0.473 |
| coverage_sampled | **0** | 0.044 | **0** | **0** | **0** | **0** |
| coverage_$\mathcal{L}$_train | 0.431 | **0.001** | **0** | **0** | **0** | 0.127 |
| coverage_$\mathcal{L}$_sampled | **0** | 0.035 | 0.004 | **0** | 0.059 | **0** |
| precision_train | 1 | **0** | **0** | **0** | **0** | **0** |
| precision_sampled | 0.025 | **0** | 0.623 | 0.104 | 0.042 | 0.004 |

Table 8: Statistical analysis of RQ 2. Pairwise comparison of using training data vs. sampled data for $\bar{\mathbf{X}}$. Each value corresponds to the p-value obtained for the Wilcoxon rank-sum test with $H_0$ that the performance of methods using training data is better than the performance of methods using sampled data. Cells printed in bold font correspond to p-values that are lower than $\alpha = 0.05/30$ (Bonferroni-adjustment) and indicate a preference towards using sampled data. Comparisons were only conducted for the methods run without post-processing.

| measure | overall | MaxBox | PRIM | Anchors | MAIRE |
|---|---|---|---|---|---|
| coverage_train | 1.00 | 1 | 0.999 | 0.445 | 0.744 |
| coverage_sampled | 0.00 | **0** | 0.987 | 0.402 | 0.003 |
| coverage_$\mathcal{L}$_train | 1.00 | 1 | 1 | 0.781 | 0.896 |
| coverage_$\mathcal{L}$_sampled | 0.00 | **0** | 0.236 | 0.476 | 0.172 |
| precision_train | 1.00 | 0.995 | 0.998 | 0.782 | 0.993 |
| precision_sampled | 0.00 | 0.011 | **0** | **0.001** | 0.381 |

Table 9: Statistical analysis of RQ 3. Pairwise comparison of using no post-processing vs. using post-processing. Each value corresponds to the p-value obtained for the Wilcoxon rank-sum test with $H_0$ that the performance of methods using no post-processing is better than the performance of methods using post-processing. Cells printed in bold font correspond to p-values that are lower than $\alpha = 0.05/60$ (Bonferroni-adjustment) and indicate a preference towards post-processing.

| method | coverage__train | coverage__sampled | coverage__$\mathcal{L}$_train | coverage__$\mathcal{L}$_sampled | precision__train | precision__sampled |
|---|---|---|---|---|---|---|
| **traindata** | 0.95 | **0** | 0.369 | **0** | **0** | **0** |
| MaxBox | 0.97 | **0** | 0.982 | **0** | 0.995 | 0.003 |
| PRIM | 1.00 | 1 | 1 | 0.452 | 0.999 | **0** |
| Anchors | 0.92 | 0.001 | 0.065 | 0.054 | **0** | **0** |
| MAIRE | 0.10 | **0** | 0.003 | **0** | **0** | 0.001 |
| **sampled** | 0.12 | **0** | **0** | **0** | **0** | **0** |
| MaxBox | 0.00 | **0** | **0** | **0** | 0.085 | 0.021 |
| PRIM | 0.45 | 0.19 | 0.262 | 0.468 | **0** | 0.001 |
| Anchors | 0.92 | **0** | 0.035 | 0.061 | **0** | **0** |
| MAIRE | 0.18 | **0** | 0.003 | **0** | **0** | 0.009 |