# FEATURE RESPONSIVENESS SCORES: MODEL-AGNOSTIC EXPLANATIONS FOR RECOURSE

**Harry Cheon**
UC San Diego

**Anneke Wernerfelt**
Haverford College

**Sorelle A. Friedler**
Haverford College

**Berk Ustun**
UC San Diego

## ABSTRACT

Consumer protection rules require companies that deploy models to automate decisions in high-stakes settings to explain predictions to decision subjects. These rules are motivated, in part, by the belief that explanations can promote *recourse* by revealing information that decision subjects can use to contest or overturn their predictions. In practice, companies provide individuals with a list of principal reasons based on feature importance derived from methods like SHAP and LIME. In this work, we show how common practices can fail to provide recourse and propose to highlight features based on their *responsiveness*—the probability that a decision subject can attain a target prediction through an arbitrary intervention on the feature. We develop efficient methods to compute responsiveness scores for any model and actionability constraints. We show that standard practices in lending can undermine decision subjects by highlighting unresponsive features and explaining predictions that are fixed.

## 1 INTRODUCTION

Machine learning models routinely automate and support decisions in consumer finance [31], employment [10, 51], and public services [65, 19, 27]. In these domains, companies are increasingly required to provide explanations to decision subjects who receive adverse outcomes (e.g., denied a loan) [1, 63, 58, 20]. In the European Union, for example, Article 86 of the AI Act [20] grants individuals a *right to explanation* in "high risk" domains [see Annex III of 20]. In the United States, the *adverse action* provision in the Equal Credit Opportunity Act mandates that lenders provide a list of "principal reasons" to consumers who are denied credit [1].

Explanations are a cornerstone of consumer protection in such domains because they may reveal information that consumers could use to exercise their broader rights [17]. In the European Union, for instance, the right to an explanation in the GDPR is meant to reveal information that consumers could use to contest their decisions or request human review [36]. Likewise, in the United States, adverse action notices are meant to support: *anti-discrimination*, by revealing that a prediction was based on protected characteristics; *rectification*, by revealing that a prediction was based on erroneous information; and *recourse*, by revealing how to attain a desired prediction in the future [57, 54].

Explainability mandates provide companies with substantial leeway on how they build explanations. In practice, companies resort to the path of least resistance, using popular feature attribution methods like SHAP and LIME to report features in *feature-highlighting explanations* for decision subjects. However, we do not yet know to what extent standard practices achieve the goals of explainability mandates. This information is necessary to guide efforts in enacting and enforcing explainability mandates, especially considering many are in early stages of development.

We study how explanations can effectively achieve one of their goals—helping consumers attain *recourse*. Our main contributions include:

1. We identify feature attribution methods can provide *reasons without recourse*—reporting "important" features that do not facilitate recourse.

2. We introduce an approach to highlight features that lead to recourse by measuring *responsiveness*— the probability that an individual can attain a target outcome by intervening on a specific feature.
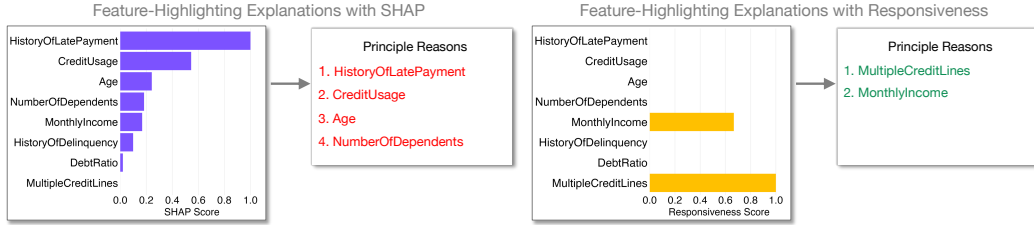
---

This is an extended version of the paper accepted at ICLR 2025.

**Figure 1:** Feature-highlighting explanations for a person denied credit by an XGBoost model on the `givemecredit` dataset in Section 5. We show explanations that highlight up to 4 features with the largest SHAP scores (left) and responsiveness scores (right). As shown, an explanation built with SHAP highlights features that the person cannot change (e.g., `Age`, `HistoryOfLatePayment`, `NumberOfDependents`) or *unresponsive* (`CreditUsage`, which can be changed but would not lead to a target prediction). In contrast, an explanation built with responsiveness scores highlights *the only* 2 features lead to a desired prediction: `MonthlyIncome` and `MultipleCreditLines`.

3. We develop methods to compute feature responsiveness scores for any classification model. Our methods can enforce complex actionability constraints that allow practitioners to control the set of interventions and their downstream effects.

4. We conduct an extensive empirical study on feature-highlighting explanations in lending. Our results show that standard methods can harm consumers by highlighting immutable and unresponsive features, and that our approach promotes recourse and transparency by highlighting responsive features and flagging predictions that are difficult or impossible to change.

5. We include a Python library to compute feature responsiveness scores, available on GitHub and installable by `pip install reachml`.

**Related Work**  Our work is related to a stream of methods to explain individual predictions [52, 44, 45, 39]. We identify these methods can inflict harm by providing individuals with *reasons without recourse*. We view reasons without recourse as a structural limitation that affects how we operationalize explainability mandates, akin to limitations of explainability that arise due to the multiplicity of predictions [46, 62, 8], the indeterminacy of explanations [12, 43], and the potential for fairwashing [4, 56, 28].

Our goals are aligned with works in algorithmic recourse, in that we seek to provide individuals with information to overturn adverse outcomes [60, 37, 61]. Many recourse methods are designed to return an *action* that an individual could perform to attain a target prediction. In contrast, our method is designed to estimate the prevalence of actions that lead to a target prediction (see Fig. 2). We construct these estimates through algorithms that sample or enumerate a set of reachable points [42]. The resulting approach is model agnostic and can be adapted them to address practical challenges related to causality [38, 15, 25] and distributional robustness [49, 50, 59].
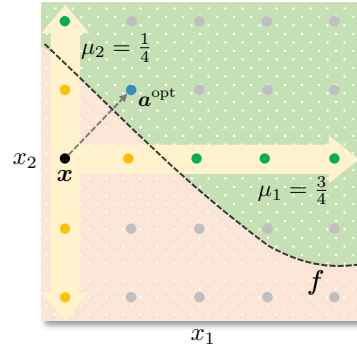


**Figure 2:** Standard methods for recourse provision return the closest action that leads to a target prediction $\boldsymbol{a}^{\mathrm{opt}}$. Our method estimates the proportion of actions on each feature that lead to a target prediction. Here, $\mu_1 = \frac{3}{4}$ and $\mu_2 = \frac{1}{4}$ because $\boldsymbol{x}$ can attain a target prediction through 3/4 actions on $x_1$, or 1/4 actions on $x_2$.

## 2 PROBLEM STATEMENT

We consider a classification task where a company uses a model $f : \mathcal{X} \to \mathcal{Y}$ to predict a label $y \in \mathcal{Y}$ from a set of *features* $\boldsymbol{x} = [x_1, x_2, \ldots, x_d] \in \mathcal{X} \subseteq \mathbb{R}^d$. We focus on tasks where each instance represents a person, and their features encode semantically meaningful characteristics. In such settings, we can assume that features are bounded. In practice, many features are bounded by definition—e.g., a binary feature such as `recent_payment` $\in \{0, 1\}$. In other cases, we can set loose bounds that apply to all decision subjects—e.g., `age` $\in [0, 120]$.

We specify the subset of individuals who are entitled to explanations in terms of a *target prediction* $\hat{y}^{\text{target}} \in \mathcal{Y}$. We assume that $\hat{y}^{\text{target}}$ represents a desirable outcome, e.g., $f(\boldsymbol{x}) = \hat{y}^{\text{target}} = 1$ if a person with features $\boldsymbol{x}$ will repay their loan. Under these conventions, companies provide an explanation to any person with features $\boldsymbol{x}$ such that $f(\boldsymbol{x}) \neq \hat{y}^{\text{target}}$. Informally, such an explanation would lead to *recourse* [60] when it contains specific information for this person to overturn an adverse outcome—e.g., by describing how to change their features to attain a point $\boldsymbol{x}'$ such that $f(\boldsymbol{x}') = \hat{y}^{\text{target}}$.

**Feature-Highlighting Explanations**   Companies comply with explainability mandates by building a *feature-highlighting explanation*—i.e., a list that contains the most important features for a specific prediction [5]. In practice, they derive the importance of each feature from post-hoc explainability methods such as LIME or SHAP. In what follows, we refer to this procedure as a *feature attribution method* and define it below.

**Definition 1.** Given a model $f : \mathcal{X} \to \mathcal{Y}$ and a point $\boldsymbol{x} \in \mathcal{X}$, a *feature attribution method* is a function $\boldsymbol{\phi} : \mathcal{X} \to \mathbb{R}^d$ that returns a vector of *feature importance scores* $\boldsymbol{\phi}(\boldsymbol{x}; f) := [\phi_1(\boldsymbol{x}; f), \dots, \phi_d(\boldsymbol{x}; f)]$. The score for each feature $\phi_j(\boldsymbol{x}; f)$ reflects its relative contribution towards the prediction. In what follows, we write $\boldsymbol{\phi}(\boldsymbol{x})$ instead of $\boldsymbol{\phi}(\boldsymbol{x}; f)$ when $f$ is clear from context.

We use the function $\boldsymbol{\phi} : \mathcal{X} \to \mathbb{R}^d$ to represent common approaches to extract feature importance scores from local explanations:

- *Local Surrogates* [see e.g., 52, 67, 66], which explain the prediction of a model $f$ at a point $\boldsymbol{x}$ by fitting a surrogate model to approximate the decision boundary of $f$ near $\boldsymbol{x}$. Given the surrogate model, we use its parameters to determine the importance scores for each feature: $\phi_j(\boldsymbol{x})$.
- *Shapley Values* [see e.g., 44, 32, 24], which cast the features of a model $f$ as "players" in a cooperative game. Each score $\phi_j(\boldsymbol{x})$ reflects the marginal contribution of feature $j$ towards the prediction $f(\boldsymbol{x})$.

Scores from these methods indicate relative importance due to the following properties:

- *Relevance*: Changing a feature $j$ with $\phi_j(\boldsymbol{x}) = 0$ does not affect the model (i.e., the feature can be dropped) [see e.g., the missingness axiom in 44].
- *Strength*: Given two features $j, k \in [d]$ such that $|\phi_j(\boldsymbol{x})| > |\phi_k(\boldsymbol{x})|$, feature $j$ has a stronger contribution to the prediction than feature $k$ [see e.g., 48].

Given a list of top-scoring features, companies convert the list into a natural language explanation [e.g., a reason code 21, 14]. In doing so, they can claim that they have met regulatory requirements by providing tailored and accessible explanations to decision subjects.

**Reasons without Recourse**   One of the limitations of feature-highlighting explanations based on importance scores is that they can highlight features that do not provide recourse. We refer to this phenomenon as *reasons without recourse*. Feature attribution methods can provide reasons without recourse due to two key blind spots:

- *Ignorance of Counterfactual Behavior*: They can assign high scores to features that are not responsive—i.e., changing the feature does not change the prediction (see e.g., Bilodeau et al. [7]).
- *Ignorance of Actionability*: They do not account for how individuals can change their features. This can lead them to assign high scores immutable features.

In practice, these failure modes can render explainability mandates counterproductive. Explanations can highlight the "wrong" features if there are other features that the decision subject can change to attain recourse. In the worst case, the explanation can be providing reasons for a fixed prediction—i.e., $f(\boldsymbol{x})$ remains the same under all feasible actions (see Table 1).

| Feature Values | | Label Counts | | Best Predictions |
|---|---|---|---|---|
| age $\geq 60$ | savings $\geq 50$K | $n_0$ | $n_1$ | $f(x_1, x_2)$ |
| 0 | 0 | 40 | 10 | 0 |
| 0 | 1 | 10 | 30 | 1 |
| 1 | 0 | 20 | 10 | 0 |
| 1 | 1 | 30 | 10 | 0 |

**Table 1:** Stylized classification task where the most accurate model assigns *fixed predictions* due to an immutable feature age $\geq 60$. We train a model to predict $y =$ repayment $(x_1, x_2) =$ (age $\geq 60$, savings $\geq 50$K) on a dataset with $n_0$ negative labels and $n_1$ positive labels. Here, individuals with age $\geq 60 = 1$ are assigned a fixed prediction, as $f(x_1, x_2) = 0$ for all reachable points $\{(1, 0), (1, 1)\}$.

## 3  MEASURING FEATURE RESPONSIVENESS

Our goal is to measure the *responsiveness* of features for decision subjects—i.e., how often the prediction of a model changes after they intervene on a given feature. However, individuals can only intervene on features in certain ways, and these changes may also have downstream effects on other features. For instance, changing married from $0 \rightarrow 1$ will change single from $1 \rightarrow 0$. Or increasing years_account_history will also cause a commensurate increase in age. In light of these challenges, we describe how decision subjects can intervene on individual features through an *intervention model*:

**Definition 2.** Given a point $\boldsymbol{x}$, we assume that individual who intervene feature $j$ will move to a new point $\boldsymbol{x}'$ where:

$$\boldsymbol{x}' = \boldsymbol{x} + \boldsymbol{a} + \boldsymbol{\varepsilon}(\boldsymbol{x}, \boldsymbol{a})$$

Here:

- $\boldsymbol{a} = [a_1, \dots, a_d] \in A_j(\boldsymbol{x})$ is an *action* that represents the deterministic components of the intervention; this includes the change in feature $j$ and its deterministic downstream effects. We assume that $a_j \neq 0$ and refer to the set of all possible actions as the *action set* $A_j(\boldsymbol{x})$.
- $\boldsymbol{\varepsilon}(\boldsymbol{x}, \boldsymbol{a})$ is a sample of the random variable that represents the stochastic component of the intervention. The sample $\boldsymbol{\varepsilon}(\boldsymbol{x}, \boldsymbol{a})$ is drawn from the probability distribution $P_{\boldsymbol{x}, \boldsymbol{a}}$.

For a fixed action $\boldsymbol{a}$, $\boldsymbol{x}'$ is a realization of a random variable. In what follows, we denote the modified features associated with a specific action $\boldsymbol{a}$ as the random variable $X_j^{\text{reach}}(\boldsymbol{x}, \boldsymbol{a})$. When the action itself is also random, we denote the resulting random variable as $X_j^{\text{reach}}(\boldsymbol{x})$.

**On Specification**  This model allows practitioners to specify feasible interventions and their deterministic and stochastic effects. They can define feasible interventions and deterministic downstream effects in the action set $A_j(\boldsymbol{x})$ through *actionability constraints*. These include *separable constraints* that only pertain to one feature (e.g., bounds and monotonicity) and *joint constraints* across multiple features (e.g., deterministic downstream effects). As shown in Table 2, we can elicit these constraints from human experts in natural language and convert them into equations that we can embed into optimization problems to enforce actionability [e.g., to search for recourse actions 60, 42].

Practitioners can define stochastic effects through the conditional probability distribution $P_{\boldsymbol{x}, \boldsymbol{a}}$. This distribution can represent probabilistic causal effects of interventions. For example, we can model the impact of employment on health insurance with $\varepsilon_{\text{has\_insurance}} \sim a_{\text{employed}} \cdot \text{Bernoulli}(\lambda)$, where whether one has health insurance largely follows their employment status. Similarly, we can model random fluctuations in features that occur between successive predictions. For instance, we can model the fluctuation in the number of bank transactions per month with $\varepsilon_{\text{n\_transactions}} \sim \text{Pois}(\lambda)$. In both cases, $P_{\boldsymbol{x}, \boldsymbol{a}}$ denotes the corresponding probability mass functions.

**Measuring Responsiveness**  Given our intervention model, we wish to score each feature by the probability that an individual attains a target prediction after performing an arbitrary intervention.

| Requirement | Example | Features | Actionability Constraint |
|---|---|---|---|
| Immutability | age cannot change | $x_j = $ age | $v_j = 0$ |
| Monotonicity | recent_payment can only increase | $x_j = $ recent_payment | $v_j \geq 0$ |
| Integrality | late_payments must be positive integer $\leq 12$ | $x_j = $ late_payments | $v_j \in \mathbb{Z}^+ \cap [0 - x_j, 12 - x_j]$ |
| Encoding Validity | preserve one-hot encoding of categorical feature housing $\in \{\text{own}, \text{rent}, \text{other}\}$ | $x_k = \mathbb{1}[\text{housing=own}]$ $x_l = \mathbb{1}[\text{housing=rent}]$ $x_m = \mathbb{1}[\text{housing=other}]$ | $v_j + x_j \in \{0, 1\}$ for $j \in \{k, l, m\}$ $\sum_{j \in \{k,l,m\}} v_j + x_j = 1$ |
| Logical Implication | if has_savings_account $=$ TRUE then savings_balance $\geq 0$ else savings_balance $= 0$ | $x_j = $ has_savings_account $x_k = $ savings_balance | $v_j + x_j \in \{0, 1\}$ $v_k + x_k \in [0, 10^{12}]$ $v_k + x_k \leq 10^{12}(x_j + v_j)$ |
| Causal Implication | if years_of_account_history increases then age will increase commensurately | $x_j = $ years_of_account_history $x_k = $ age | $x_j + v_j \leq x_k + \delta_k$ $\delta_k \in [0, 100]$ |

**Table 2:** Examples of actionability constraints on semantically meaningful features for a lending task. Each constraint can be expressed in natural language and embedded into an optimization problem using standard techniques in mathematical programming [see, e.g., 64]. See Section B for more examples.

**Definition 3.** Given a model $f : \mathcal{X} \to \mathcal{Y}$, a point $\boldsymbol{x} \in \mathcal{X}$, a feature $j \in [d]$, its action set $A_j(\boldsymbol{x})$ and the downstream distribution $P_{\boldsymbol{x},\boldsymbol{a}}$, the *responsiveness score* of feature $j$ measures the probability that an intervention on feature $j$ attains the target prediction:

$$\mu_j(\boldsymbol{x}) := \Pr\Big(f\big(X_j^{\text{reach}}(\boldsymbol{x}, \boldsymbol{a})\big) = \hat{y}^{\text{target}} \mid \boldsymbol{a} \in A_j(\boldsymbol{x})\Big)$$

Here, a score of $\mu_j(\boldsymbol{x}) = 0$ means that changing feature $j$ cannot achieve $\hat{y}^{\text{target}}$, while $\mu_j(\boldsymbol{x}) = 1$ means any intervention on $j$ will achieve $\hat{y}^{\text{target}}$.

**Benefit for Consumer Protection**   When we construct feature-highlighting explanations using the top-$k$ most responsive features, we reveal the $k$ most promising paths to recourse. By construction, they are features where arbitrary interventions are most likely to lead to a target prediction. We make no assumptions on how individuals will change their features beyond actionability constraints specified in $A_j(\boldsymbol{x})$. This is because modeling how individuals will intervene on features is not feasible; it cannot be verified a priori.

Contrary to existing methods, our approach provides explanations *only* to individuals with recourse— i.e., we would never provide reasons without recourse. In effect, we can detect instances where providing feature-highlighting explanations can be misleading or harmful by checking if the responsiveness score $\mu_j(\boldsymbol{x}) = 0$ for all features.

**Remark 1.** *Given a model $f : \mathcal{X} \to \mathcal{Y}$, denote its feature responsiveness scores at the point $\boldsymbol{x} \in \mathcal{X}$ as $\mu_1(\boldsymbol{x}) \dots \mu_d(\boldsymbol{x})$. If $\mu_j(\boldsymbol{x}) = 0$ for all $j \in [d]$, then either:*

*(a) $f$ assigns a fixed prediction to $\boldsymbol{x}$, or*

*(b) $f$ can only provide recourse to $\boldsymbol{x}$ through an intervention on two or more features.*

According to Remark 1, there are two scenarios where $\mu_j(\boldsymbol{x}) = 0$ for all $j \in [d]$. We can employ different strategies to mitigate harm in each case. In Case (a), where individuals receive fixed predictions, we can withhold explanations and notify developers or regulators. In Case (b), where individuals can only overturn their prediction by intervening on multiple features at the same time, we can include a warning against assuming feature independence.

The responsiveness score $\mu_j(\boldsymbol{x})$ depends on the actionability constraints that characterize $A_j(\boldsymbol{x})$— i.e., responsiveness scores can change under a different set of constraints. Hence, if constraints are misspecified (e.g., ignoring monotonicity constraints), the scores can lead to misleading conclusions. In tasks where downstream effects are deterministic, we can mitigate this effect by encoding *indisputable constraints* based on feature encodings or physical limits. Then, the corresponding $\mu_j(\boldsymbol{x})$ represents an upper bound on the true responsiveness of feature $j$—i.e., decision subjects flagged with a fixed prediction will also have a fixed prediction under more stringent constraints.

## 4   COMPUTING SCORES WITH REACHABLE SETS

We now introduce an approach to compute responsiveness scores for any model. We compute the responsiveness score of feature $j$ using its reachable set $R_j(\boldsymbol{x})$—the set of reachable points through interventions on $j$ (see Fig. 3). We can generate the reachable set $R_j(\boldsymbol{x})$ either by enumerating all possible points—when features are discrete and interventions do not have stochastic effects—or sampling. Given $R_j(\boldsymbol{x})$, we can compute the responsiveness score as:

$$\mu_j(\boldsymbol{x}) := \mathbb{E}_{\boldsymbol{x}' \sim X_j^{\text{reach}}(\boldsymbol{x})}\big[\mathbb{1}[f(\boldsymbol{x}') = \hat{y}^{\text{target}}]\big] \tag{1}$$

This approach has several benefits. It is model agnostic; given the reachable sets for each feature, computing the reachable set only requires query access to the model. We can also amortize the cost of generating reachable points by generating the reachable sets *once* and (re)using it to compute responsiveness scores for any model (e.g., during model selection).

In practice, enumerating or sampling reachable points can be challenging. We often need to consider points in regions with little structure—points may have both discrete and continuous dimensions and obey non-convex constraints. Furthermore, when we have stochastic downstream effects, we must check that both interventions and their downstream effects are feasible under actionability constraints.

**Figure 3:** Stylized example showing how to compute responsiveness scores for a classification model with three features n_loans, guarantor and age. The reachable set $R_j(\boldsymbol{x})$ all points that can be attained from $\boldsymbol{x} = (3, 0, 24)$ by intervening on feature $j$, and $R_3(\boldsymbol{x}) = \varnothing$ because age is immutable. Given a model $f$, we compute the responsiveness score of each feature by querying its predictions over points in their reachable set $R_j(\boldsymbol{x})$.

We overcome these challenges by casting the generation of reachable points as repeated optimization problems.

**Sampling** We present a procedure to sample reachable points in Algorithm 1. Given a point $\boldsymbol{x}$ and an action set $A_j(\boldsymbol{x})$, this procedure returns a uniform sample of $N$ reachable points via rejection sampling. In Section 4, it calls the $\mathsf{Sample1DAction}(\boldsymbol{x}, A_j)$ routine to propose a candidate deterministic change $\boldsymbol{a}$ that obeys separable constraints such as bounds and integrality. We then sample its stochastic downstream effect $\boldsymbol{\varepsilon}$ in Section 4. In Section 4, it then calls the $\mathsf{CheckFeasibility}$ routine to check if both the intervention and the downstream effect obey actionability constraints by solving a mixed-integer program. The procedure terminates once it has sampled $N$ reachable points through interventions on $j$. Given $R_j(\boldsymbol{x})$, we can recover an unbiased estimate of the responsiveness score for feature $j$ and a model $f$ as $\hat{\mu}_j(\boldsymbol{x}) := \frac{1}{N}\sum_{\boldsymbol{x}' \in \hat{R}_{\boldsymbol{x}}(j)} \mathbb{1}[f(\boldsymbol{x}') = \hat{y}^{\text{target}}]$. We can set the sample size $N$ to ensure practical guarantees on how reliably we flag fixed predictions (Remark 1) as described in Section A.2.

---

**Algorithm 1** Sample Reachable Points

**Require:** $\boldsymbol{x} \in \mathcal{X}$          *point*
**Require:** $A_j(\boldsymbol{x})$      *action set for feature $j$*
**Require:** $P_{\boldsymbol{x},\boldsymbol{a}}$   *stochastic downstream effect dist. for $j$*
**Require:** $N \in \mathbb{N}$     *sample size (see Section A.2)*

    $\hat{R}_j \leftarrow \varnothing$
1: **repeat**
2:     $\boldsymbol{a} \leftarrow \mathsf{Sample1DAction}(\boldsymbol{x}, A_j)$
3:     $\boldsymbol{\varepsilon} \leftarrow \mathsf{SampleDownstream}(\boldsymbol{x}, \boldsymbol{a}, P_{\boldsymbol{x},\boldsymbol{a}})$
4:     **if** $\mathsf{CheckFeasibility}(\boldsymbol{x}, \boldsymbol{a} + \boldsymbol{\varepsilon}, A_j)$ **then**
5:         $\hat{R}_j \leftarrow \hat{R}_j \cup \{\boldsymbol{x} + \boldsymbol{a} + \boldsymbol{\varepsilon}\}$
6:     **end if**
7: **until** $|\hat{R}_j| = N$
**Output** $\hat{R}_j$    *N reachable points via actions on $j$*

---

**Enumeration** We present a procedure to enumerate a reachable set $R_j(\boldsymbol{x})$ in Algorithm 2 for discrete features with deterministic downstream effects. Given the action set, which encodes all actionability constraints (including deterministic downstream effects), the procedure enumerates reachable points for feature $j$ by repeatedly solving the following discrete optimization problem:

$$\mathsf{Find1DAction}(\boldsymbol{x}, A_j) := \underset{\boldsymbol{a} \in A_j(\boldsymbol{x})}{\operatorname{argmin}} \|\boldsymbol{a}\|_1$$

---

**Algorithm 2** Enumerate Reachable Points

**Require:** $\boldsymbol{x} \in \mathcal{X}$         *point*
**Require:** $A_j(\boldsymbol{x})$   *action set for discrete feature $j$*
    $R_j \leftarrow \varnothing, A_j \leftarrow A_j(\boldsymbol{x})$
1: **repeat**
2:     $\boldsymbol{a}^* \leftarrow \mathsf{Find1DAction}(\boldsymbol{x}, A_j)$
3:     $R_j \leftarrow R_j \cup \{\boldsymbol{x} + \boldsymbol{a}^*\}$
4:     $A_j \leftarrow A_j \setminus \{\boldsymbol{a}^*\}$
5: **until** $\mathsf{Find1DAction}(\boldsymbol{x}, A_j)$ is infeasible
**Output** $R_j$   *all reachable points via actions in $j$*

---

We formulate $\mathsf{Find1DAction}(\boldsymbol{x}, A_j)$ as a mixed-integer program, and update it at each iteration with a "no good" constraint to remove previous optima in Section 4 (see Section A.1 for exact formulation). We use each action to add a reachable point to $R_j(\boldsymbol{x})$ and use the final set to calculate *exact* responsiveness scores. We adapt a method to enumerate the reachable set for all features from Kothari et al. [42], but is more tractable as we only enumerate points that can we can attain through interventions on each feature.

**Extensions** One of the benefits of reachable sets is that we easily customize scores to meet additional requirements. One such requirement is *monotonicity*, i.e., if a person is guaranteed a target prediction by increasing (or decreasing) a feature beyond a threshold value. In the simplest case, we can account

for properties through operations like filtering or weighing (see e.g., Section 6). In general, we can construct responsiveness scores that address practical challenges given additional inputs:

- *Individual Preferences*: Given a cost function that captures the difficulty of actions in each direction, we can highlight features that are easier to change (i.e., least costly $k$ features) using a cost-weighed score: $\mu_j^{\text{cost}}(\boldsymbol{x}; \text{ cost}) = \sum_{\boldsymbol{x}' \in R_j(\boldsymbol{x})} \text{cost}(\boldsymbol{x}'; \boldsymbol{x}) \cdot \mathbb{1}[f(\boldsymbol{x}') = \hat{y}^{\text{target}}]$.
- *Distributional Robustness*: Given a general reachable set $R(\boldsymbol{x})$ that contains all points that we could reach through interventions on any feature, we highlight features that attain a target prediction regardless of how other features change through the *robust score*: $\mu_j^{\text{robust}}(\boldsymbol{x}) = \min_{\boldsymbol{\delta} \in \Delta_{-j}} \mathbb{E}_{X' \sim R_j(\boldsymbol{x})}[\mathbb{1}[f(X' + \boldsymbol{\delta}) = \hat{y}^{\text{target}}]]$, where $\Delta_{-j} := \{\boldsymbol{\delta} \in \mathbb{R}^d \mid \delta_j = 0, \|\boldsymbol{\delta}\| < \varepsilon\}$.

## 5 EXPERIMENTS

We present an empirical study on the responsiveness of explanations. Our results reveal the limitations of existing feature attribution methods and show how our approach can support recourse and flag fixed predictions. We include details in Section B, and code to reproduce our results on GitHub.

**Setup** We work with three publicly available classification datasets from consumer finance. Here, each instance represents a consumer and the label indicates if they will repay a loan. We consider discrete version of each dataset in which we can compute exact responsiveness scores and certify if each person has recourse. Given these datasets, we define *inherent actionability constraints* which reflect indisputable requirements that apply to all individuals (e.g., no changes to immutable attributes, preserve feature encoding, and adhere to deterministic causal effects).

We split each dataset into a training sample (80%; to train models) and a test sample (20%; to evaluate out-of-sample performance). We fit models using (1) *logistic regression* (LR), (2) *XGBoost* (XGB), and (3) *random forests* (RF). For each model, we construct a feature-highlighting explanation for each person who is denied credit in the dataset that includes up to *four features*; if all features have a score of 0, we do not present an explanation for that individual. The choice of *up to four* features reflects the recommended number of reasons to show in an adverse action notice by the U.S. Consumer Finance Protection Bureau [see 2]. We include the top-4 scoring features from the following methods:

- *Feature Responsiveness* (RESP): We compute responsive scores from complete reachable sets that we enumerate using Algorithm 2.
- *Standard Feature Attribution*: We consider model-agnostic methods that are widely used in the lending industry [23]: SHAP [44]; and LIME [52].
- *Actionable Feature Attribution*: We consider *action-aware* variants of SHAP and LIME: SHAP-AW and LIME-AW. They aim to highlight responsive features by $\phi_j(\boldsymbol{x}) \leftarrow 0$ for immutable features.

| Dataset | Metrics | LR | RF | XGB |
|---|---|---|---|---|
| heloc | % Denied | 56.1% | 58.3% | 57.0% |
| $n = 5,842$ | ↳ % No Recourse | 22.2% | 31.3% | 53.1% |
| $d = 43$ | ↳ % 1-D Recourse | 41.0% | 31.7% | 25.3% |
| FICO [22] | ↳ % $n$-D Recourse | 36.8% | 37.0% | 21.6% |
| german | % Denied | 22.9% | 17.5% | 22.0% |
| $n = 1,000$ | ↳ % No Recourse | 7.4% | 28.6% | 11.8% |
| $d = 36$ | ↳ % 1-D Recourse | 74.2% | 48.0% | 68.2% |
| Dua & Graff [16] | ↳ % $n$-D Recourse | 18.3% | 23.4% | 20.0% |
| givemecredit | % Denied | 24.6% | 24.7% | 24.8% |
| $n = 120,268$ | ↳ % No Recourse | 15.6% | 0.2% | 11.5% |
| $d = 23$ | ↳ % 1-D Recourse | 72.4% | 93.2% | 76.0% |
| Kaggle [35] | ↳ % $n$-D Recourse | 12.0% | 6.6% | 12.5% |

**Table 3:** Overview of paths to recourse for individuals who would receive an explanation for each dataset and model. We report % *Denied*, % of denied individuals; *% No Recourse*, % of denied with a fixed prediction (i.e., who have no recourse under any explanation); *% 1-D*, % of denied individuals who can overturn their prediction by changing 1 feature (i.e., who could have recourse from a feature-highlighting explanation); and *% n-D*, % of denied individuals who can only overturn their prediction by changing 2 or more features simultaneously.

**On the Limits of Feature-Highlighting Explanations** Our results in Table 3 highlight how current practices to comply with explainability mandates can help consumers achieve recourse. As shown, there is no case— i.e., for any model, any dataset, and any explanation method—where all individuals that receive feature-highlighting explanations could attain the target prediction through a single-feature intervention. Some require joint interventions. Others have no path to recourse.

On the heloc dataset, for example, a lender who uses an LR model would provide feature-highlighting explanations to 56.1% of applicants. Among these individuals, 41.0% could attain a desired prediction by changing single feature, 36.8% could only do so by changing 2 or more features simultaneously, and the model assigns a fixed prediction to the remaining 22.2%.

| Dataset | Metrics | LR | | | | | XGB | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | All Features | | Actionable Features | | | All Features | | Actionable Features | | |
| | | LIME | SHAP | LIME-AW | SHAP-AW | RESP | LIME | SHAP | LIME-AW | SHAP-AW | RESP |
| heloc | % Presented with Explanations | 100.0% | 100.0% | 100.0% | 100.0% | 41.0% | 100.0% | 100.0% | 100.0% | 100.0% | 25.3% |
| $n = 5,842$ | ↳ % All Features Unresponsive | 92.7% | 77.3% | 76.8% | 70.3% | 0.0% | 93.2% | 82.3% | 80.0% | 79.6% | 0.0% |
| $d = 43$ features | ↳ % At Least 1 Feature Responsive | 7.3% | 22.7% | 23.2% | 29.7% | 100.0% | 6.8% | 17.7% | 20.0% | 20.4% | 100.0% |
| $d_A = 31$ mutable | ↳ % All Features Responsive | 0.0% | 0.0% | 0.0% | 0.2% | **100.0%** | 0.0% | 0.0% | 0.0% | 0.0% | **100.0%** |
| FICO [22] | ↳ # Features Highlighted | 4.0 | 4.0 | 4.0 | 4.0 | 2.3 | 4.0 | 4.0 | 4.0 | 4.0 | 2.5 |
| german | % Presented with Explanations | 100.0% | 100.0% | 100.0% | 100.0% | 74.2% | 100.0% | 100.0% | 100.0% | 100.0% | 68.2% |
| $n = 1,000$ | ↳ % All Features Unresponsive | 91.7% | 100.0% | 59.4% | 65.1% | 0.0% | 100.0% | 99.1% | 70.5% | 67.3% | 0.0% |
| $d = 36$ features | ↳ % At Least 1 Feature Responsive | 8.3% | 0.0% | 40.6% | 34.9% | 100.0% | 0.0% | 0.9% | 29.5% | 32.7% | 100.0% |
| $d_A = 9$ mutable | ↳ % All Features Responsive | 0.0% | 0.0% | 0.0% | 0.0% | **100.0%** | 0.0% | 0.0% | 0.0% | 0.0% | **100.0%** |
| Dua & Graff [16] | ↳ # Features Highlighted | 4.0 | 4.0 | 4.0 | 4.0 | 1.8 | 4.0 | 4.0 | 4.0 | 4.0 | 1.8 |
| givemecredit | % Presented with Explanations | 100.0% | 100.0% | 100.0% | 100.0% | 72.4% | 100.0% | 100.0% | 100.0% | 100.0% | 76.0% |
| $n = 120,268$ | ↳ % All Features Unresponsive | 65.5% | 46.8% | 33.1% | 56.0% | 0.0% | 41.8% | 43.3% | 31.6% | 30.6% | 0.0% |
| $d = 23$ features | ↳ % At Least 1 Feature Responsive | 34.5% | 53.2% | 44.0% | 66.9% | 100.0% | 58.2% | 56.7% | 68.4% | 69.4% | 100.0% |
| $d_A = 13$ mutable | ↳ % All Features Responsive | 0.0% | 0.0% | 0.0% | 22.8% | **100.0%** | 0.0% | 0.0% | 4.2% | 13.2% | **100.0%** |
| Kaggle [35] | ↳ # Features Highlighted | 4.0 | 4.0 | 4.0 | 4.0 | 2.4 | 4.0 | 4.0 | 4.0 | 4.0 | 2.6 |

**Table 4:** Responsiveness of feature-highlighting explanations for LR and XGB models for all methods and datasets. We defer results for RF to Section C.1 for clarity. For each model, we generate explanations that highlight up to 4 top-scoring features from a given method. We report the proportion of individuals receiving an explanation (*% Presented with Explanations*) and the mean number of features in each explanation (*# Features Highlighted*). We also show the proportion of instances where all features are unresponsive (*% All Features Unresponsive*) highlighting positive values; at least one feature is responsive (*% At Least 1 Feature Responsive*), or all features are responsive (*% All Features Responsive*) highlighting the **best value**.

These results reflect the *best* we can hope for when providing recourse with feature-highlighting explanations. Here, the 41.0% of individuals who could achieve recourse by a single-feature intervention can only do so if we construct explanations with an *ideal* method that assigns the highest scores to responsive features, and do not face additional actionability constraints.

**On Explanations with Feature Attribution Scores**   Our results show how standard methods for feature attribution can highlight features that are uninformative or misleading. Given the LR model on the heloc dataset, we find that 92.7% and 77.3% of explanations from LIME and SHAP fail to highlight even one responsive feature. This stems from two issues:

- *Low Scores for Responsive Features*: Under the LR model on the heloc dataset, 41.0% of denied individuals could be approved by altering a single feature. However, LIME and SHAP do not highlight these features because they assign higher scores to other features (see Section C.2).
- *Fixed Predictions*: Under the LR model on the heloc dataset, 22.2% of denied individuals cannot be approved under any intervention as they receive fixed prediction. These are instances where LIME and SHAP (and their variants) can inflict harm by highlighting mutable features. For example, one individual who is assigned a fixed prediction would receive an explanation that highlights mutable features such as AvgYearsInFile and NetFractionRevolvingBurden under SHAP, which gives the impression that intervening on them could lead to approval.

**On Adapting Existing Methods**   Seeing how feature attribution methods like LIME and SHAP can highlight features that are important but immutable, we study the potential to improve responsiveness using *action-aware* variants SHAP-AW and LIME-AW. Following a common belief that we can enforce actionability post-hoc [e.g., 47], we construct explanations using only actionable features. In Table 4, we see that SHAP-AW and LIME-AW can highlight more responsive features. Given the LR model in heloc, for example, this strategy improves the proportion of explanations that contain at least one responsive feature by 7.0% (i.e., 29.7% of SHAP-AW vs. 22.7% for SHAP). One shortcoming of this approach is that we must filter features based on their actionability for all individuals, which may overlook features that is actionable for some individuals but not others.

**On Explanations with Responsiveness Scores**   Our results show how practitioners can use our approach to comply with regulatory requirements and address the limitations of feature attribution methods. When we construct feature-highlighting explanations using responsiveness scores, we present individuals with explanations that only contain responsive features (100% on the *% All Features Responsive* metric across datasets and models in Table 4). In contrast, only 0.2% of SHAP-AW explanations of the LR model in heloc were fully responsive. For the remaining 99.8%, each

explanation contains at least one unresponsive feature that could lead individuals to intervene without achieving the target prediction.

Explanations based on responsiveness scores contain the *most* responsive features that one can change independently to achieve recourse. In effect, we only provide explanations to individuals who can achieve recourse with a single-feature intervention. This may result in explanations that highlight fewer features on average. For example, RESP explanations for the LR model on `heloc` contained an average of 2.3 (out of 4) features. This behavior can mitigate harm as we avoid presenting explanations to individuals with fixed predictions (i.e., cannot change their predictions), or to individuals who could only do so with joint interventions.

## 6   ON THE LIMITS OF FEATURE HIGHLIGHTING FOR RECOURSE

Current mandates only require explanations to contain a list of most important features. In effect, most explanations lack details on *how* individuals should intervene on them. Consider a person who receives an explanation that highlights income. In this case, most consumers would assume that increasing income would eventually lead to approval. In practice, this can yield counterproductive results—a responsive feature may be responsive in a way that is not monotonic (e.g., increase income by at least \$1,000 but not more than \$2,500) or intuitive (e.g., decrease income to be approved). As seen in this example, feature highlighting explanations can only provide recourse if features are: (1) responsive, (2) monotonic and (3) intuitive.

Building upon our results from Section 5, **we show that even when methods like LIME and SHAP highlight responsive features, the necessary interventions to obtain recourse are not immediately apparent to consumers**. These results also highlight how we can use our machinery to check for more complicated notions of responsiveness (e.g., monotonicity).

**Setup**   We use the same setup in Section 5 and fit an XGB model for a version of the `givemecredit` dataset $n = 23,459$ where we do not binarize continuous features. We construct feature highlighting explanations for each individual denied credit that contain up to four features based on scores from SHAP, LIME, and RESP. We estimate the responsiveness score (RESP) for each feature using a sample of $N = 500$ reachable points $R_j(\boldsymbol{x})$ that we generate using Algorithm 1. Our choice of $N = 500$ ensures that there is a 99% chance that a feature that we claim is unresponsive has a true responsiveness $\leq 0.01$—i.e., at most of 1% of actions could lead to recourse (see Section A.2). In addition to measuring responsiveness of each feature under arbitrary changes, we use $R_j(\boldsymbol{x})$ to evaluate whether responsiveness is monotonic or intuitive. We verify that a feature is intuitively responsive by finding at least one reachable point in the intuitive direction that attains the target prediction. We verify monotonicity of responsiveness by searching for a threshold value where all reachable points with feature $j$ less/greater than the threshold attain the target prediction. When we construct explanations with RESP through $R_j(\boldsymbol{x})$, we only include features satisfying the necessary conditions for recourse depending on the additional information we assume is provided; if none exist, we don't not construct an explanation.

**On the Need for Additional Information**   Our results reveal that omitting additional information on interventions undermines the value of feature-highlighting explanations for recourse. In Table 5, we see that there is no case where an individual can reliably achieve recourse using a feature-highlighting explanation based on SHAP or LIME. In particular, 0% of individuals receive explanations where all four feature are responsive, monotonic, and intuitive, and only 6% of individuals receive an explanation where at least 1 feature obeys these conditions (**Features Only**). This is either because interventions only lead to the target prediction when a feature takes on very specific values, or because the feature must be changed in counterintuitive ways.

*Degree of Change*: Explanations can fail to provide a reliable path to recourse when they do not include information on *how much* to change the value of responsive features. Consider credit utilization (CreditUtil), which is in 32.3% of explanations built using SHAP scores. Our analysis reveals that it is responsive in 10.7% of cases, but responsive *and monotonic* in 9.0% of cases. In other words, 1.7% of individuals can only reliably intervene on this feature to attain a desired prediction when they have additional information on how much to change their credit utilization. For instance, we point to an individual with CreditUtil $= 0.99$ who would be approved if they decrease the value of
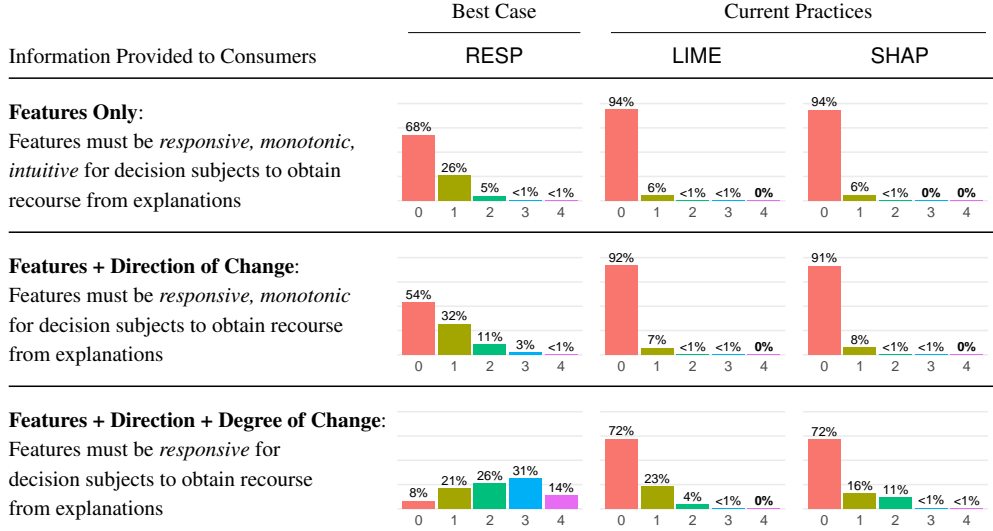
| Information Provided to Consumers | Best Case | Current Practices | |
|---|---|---|---|
| | RESP | LIME | SHAP |

**Features Only**:
Features must be *responsive, monotonic, intuitive* for decision subjects to obtain recourse from explanations

RESP: 68% (0), 26% (1), 5% (2), <1% (3), <1% (4)
LIME: 94% (0), 6% (1), <1% (2), <1% (3), 0% (4)
SHAP: 94% (0), 6% (1), <1% (2), 0% (3), 0% (4)

**Features + Direction of Change**:
Features must be *responsive, monotonic* for decision subjects to obtain recourse from explanations

RESP: 54% (0), 32% (1), 11% (2), 3% (3), <1% (4)
LIME: 92% (0), 7% (1), <1% (2), <1% (3), 0% (4)
SHAP: 91% (0), 8% (1), <1% (2), <1% (3), 0% (4)

**Features + Direction + Degree of Change**:
Features must be *responsive* for decision subjects to obtain recourse from explanations

RESP: 8% (0), 21% (1), 26% (2), 31% (3), 14% (4)
LIME: 72% (0), 23% (1), 4% (2), <1% (3), 0% (4)
SHAP: 72% (0), 16% (1), 11% (2), <1% (3), <1% (4)

**Table 5:** Distribution of features that are responsive, monotonic and intuitive in feature-highlighting explanations using LIME, SHAP, RESP. We plot the percentage of explanations given to consumers denied credit by the XGB model in the `givemecredit` dataset with $k \in \{0, \ldots, 4\}$ features that are *responsive* (can be changed to attain recourse), *monotonic* (all changes below or beyond a threshold value will lead to recourse), and *intuitive* (will lead to recourse if they are changed in a direction that aligns with common expectations). We show that 94% of consumers who receive feature-highlighting explanations built from LIME and SHAP are unlikely overturn their predictions. We can reduce this value to 72% by including additional information on the degree and direction of change. Under RESP, we only include features that meet the necessary conditions for recourse; if no features meet these conditions, we do not provide explanations for the consumer. Including additional information will decrease this proportion to 8%.

this feature to CreditUtil $\in (0.00, 0.50) \cup (0.65, 0.68)$. Without information on the degree of change, they may still be denied if they decrease their usage to the values of CreditUtil $\cup (0.50, 0.65)$.

*Direction of Change*: Explanations can also fail to provide a reliable path to recourse when they do not inlcude information on whether the decision subject should intervene on features by increasing or decreasing them. In this case, individuals who are shown responsive features may fail to obtain the target prediction because they must change features in a counterintuitive direction. Consider Income, which is in 48.1% of explanations built using SHAP and is responsive in 20.6% of all cases (i.e., 43% of explanations with Income). In general, an individual who is denied credit and shown this feature would naturally assume that they can be approved by increasing this value. Yet, our analysis reveals this is only the case for 9.9% of individuals; the remaining 10.7% of individuals, where Income is responsive, can only be approved by *decreasing* Income.

**Discussion**    When we construct explanations using methods like SHAP and LIME, we cannot reliably tell when features are responsive, monotonic, or intuitive. In contrast, with RESP, we can verify whether each feature is monotonically or intuitively responsive through its reachable set $R_j(\boldsymbol{x})$. If there are no features that meet these conditions, we would refrain from providing explanations. Hence, RESP represents "optimal conditions" for feature-highlighting explanations. However, we see that, without details on the magnitude or the direction of change, highlighting features alone is unlikely to provide recourse. Specifically, we cannot provide explanations to 68% of denied individuals (Table 5). They would not benefit from feature-highlighting explanations because none of their features are responsive, monotonic and intuitive. These results underscore the need to include additional information when explanations are meant to support recourse.

We could impose some of these conditions by enforcing constraints on how we train the model—i.e., we could ensure monotonicity by using a linear classifier like LR rather than XGB. Alternatively, we could use custom responsiveness scores to highlight features that meet all these conditions by inspecting their reachable sets $R_j(\boldsymbol{x})$, as we have done in this section. The RESP column in Table 5 shows how effective custom responsiveness scores can be as a standalone solution. When there are no features that satisfy these conditions, we could highlight features that achieve weaker forms

of responsiveness alongside additional information. In Table 5, we see that including additional information can reduce the proportion of cases we refrain from explaining to 8%. This provides an alternative approach to ensure these conditions in a way that would not interfere with model development.

## 7 CONCLUDING REMARKS

Explanations are often seen as safeguards in consumer-facing applications as they can reveal information that can help exercise their broader rights pertaining to anti-discrimination, correction of erroneous information and recourse [17].

However, our findings suggest that current mandates may be insufficient in achieving their stated goals. We showed that feature-highlighting explanations can fail to help consumers achieve a target prediction—providing *reasons without recourse*. This arises because feature attribution methods like SHAP and LIME overlook whether features are actionable and highlight immutable ones. They fail to capture counterfactual behavior and highlight unresponsive features. These explanations may also lead to harm by explaining fixed predictions. We further demonstrated that even when explanations highlight responsive features, feature-highlighting explanations may still fail to provide recourse if they omit information on how to change each feature.

Evaluating whether mandates accomplish their goals is especially important as explanation mandates are often designed to achieve objectives that can be addressed using other techniques—e.g., anti-discrimination via auditing [53, 6, 55] or searching for less discriminatory models [9, 26].

**Use Cases for Responsiveness Scores**   Our work has primarily focused on consumer finance applications because there are long-standing regulations on explanations and recourse in place. More broadly, we can draw on responsiveness scores to describe how models behave with respect to user interactions.

- *Concept Annotations*: We can use responsiveness scores to select which concepts to confirm when performing test-time interventions under a constrained budget in concept-based models [40, 34]. Interventions can prioritize confirming concepts with high responsiveness.
- *Selective Feature Reporting*: Users can decide whether to provide optional features for personalized predictions based on responsiveness scores [33]; a high responsiveness score suggests that it is beneficial for the user to report the feature.
- *Model Steering*: Users can identify features they can use to collectively steer model behavior with responsiveness scores [30].
- *Strategic Classification*: Model developers can preemptively identify features that users can manipulate to "game" their predictions and act upon them (e.g., make responsive features immutable on the platform) [29].

**Limitations**   In applications like lending, actionability of features can be ambiguous; practitioners need to make assumptions on how features can change. In this work, we measured the responsiveness of features with respect to a conservative set of assumptions—indisputable constraints on how individuals can change their features. In this regime, our responsiveness scores can flag individuals with fixed predictions; but, we may not guarantee recourse as consumers may face additional constraints that we did not enforce. In practice, we can mitigate these issues by highlighting features that exceed a minimal level of responsiveness, or by eliciting constraints from decision subjects [see e.g., 18, 13, 41].

## REFERENCES

[1] 12 cfr part 1002 - equal credit opportunity act (regulation b). https://www.consumerfinance.gov/rules-policy/regulations/1002/2/, . Accessed: 2024-07-16.

[2] Comment for 1002.9 - notifications. https://www.consumerfinance.gov/rules-policy/regulations/1002/interp-9/#9-b-1-Interp-1, . Accessed: 2024-07-16.

[3] Alan Agresti and Brent A Coull. Approximate is better than "exact" for interval estimation of binomial proportions. *The American Statistician*, 52(2):119–126, 1998.

[4] Ulrich Aïvodji, Hiromi Arai, Olivier Fortineau, Sébastien Gambs, Satoshi Hara, and Alain Tapp. Fairwashing: the risk of rationalization. In *International Conference on Machine Learning*, pp. 161–170. PMLR, 2019.

[5] Solon Barocas, Andrew D Selbst, and Manish Raghavan. The hidden assumptions behind counterfactual explanations and principal reasons. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 80–89, 2020.

[6] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilović, et al. Ai fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development*, 63(4/5):4–1, 2019.

[7] Blair Bilodeau, Natasha Jaques, Pang Wei Koh, and Been Kim. Impossibility theorems for feature attribution. *Proceedings of the National Academy of Sciences*, 121(2):e2304406120, 2024.

[8] Emily Black, Manish Raghavan, and Solon Barocas. Model multiplicity: Opportunities, concerns, and solutions. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pp. 850–863, 2022.

[9] Emily Black, John Logan Koepke, Pauline T Kim, Solon Barocas, and Mingwei Hsu. Less discriminatory algorithms. *Geo. LJ*, 113:53, 2024.

[10] Miranda Bogen and Aaron Rieke. Help wanted: An examination of hiring algorithms, equity, and bias. *Upturn, December*, 7, 2018.

[11] Lawrence D Brown, T Tony Cai, and Anirban DasGupta. Interval estimation for a binomial proportion. *Statistical science*, 16(2):101–133, 2001.

[12] Marc-Etienne Brunet, Ashton Anderson, and Richard Zemel. Implications of model indeterminacy for explanations of automated decisions. *Advances in Neural Information Processing Systems*, 35:7810–7823, 2022.

[13] Giovanni De Toni, Paolo Viappiani, Stefano Teso, Bruno Lepri, and Andrea Passerini. Personalized algorithmic recourse with preference elicitation. *arXiv preprint arXiv:2205.13743*, 2022.

[14] Louis DeNicola. What are credit score reason codes? myFico, February 2022. URL https://www.myfico.com/credit-education/blog/reason-codes.

[15] Ricardo Dominguez-Olmedo, Amir H Karimi, and Bernhard Schölkopf. On the adversarial robustness of causal algorithmic recourse. In *International Conference on Machine Learning*, pp. 5324–5342. PMLR, 2022.

[16] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.

[17] Lilian Edwards and Michael Veale. Slave to the algorithm: Why a right to an explanation is probably not the remedy you are looking for. *Duke L. & Tech. Rev.*, 16:18, 2017.

[18] Seyedehdelaram Esfahani, Giovanni De Toni, Bruno Lepri, Andrea Passerini, Katya Tentori, and Massimo Zancanaro. Exploiting preference elicitation in interactive and user-centered algorithmic recourse: An initial exploration. *arXiv preprint arXiv:2404.05270*, 2024.

[19] Virginia Eubanks. *Automating inequality: How high-tech tools profile, police, and punish the poor*. St. Martin's Press, 2018.

[20] Council of the European Union European Parliament. Regulation (eu) 2024/1689. https://eur-lex.europa.eu/eli/reg/2024/1689/oj. Accessed: 2024-08-30.

[21] Experian. Experian/fair, isaac model v2 score factor codes. https://www.trudiligence.com/docs/ficoscorefactors.pdf.

[22] FICO. Explainable machine learning challenge, 2018. URL https://community.fico.com/s/explainable-machine-learning-challenge.

[23] FinRegLab. Empirical white paper: Explainability and fairness: Insights from consumer lending. Technical report, FinRegLab, July 2023. URL https://finreglab.org/wp-content/uploads/2023/12/FinRegLab_2023-07-13_Empirical-White-Paper_Explainability-and-Fairness_Insights-from-Consumer-Lending.pdf.

[24] Fabian Fumagalli, Maximilian Muschalik, Patrick Kolpaczki, Eyke Hüllermeier, and Barbara Hammer. Shap-iq: Unified approximation of any-order shapley interactions. *Advances in Neural Information Processing Systems*, 36, 2024.

[25] Sainyam Galhotra, Romila Pradhan, and Babak Salimi. Explaining black-box algorithms using probabilistic contrastive counterfactuals. In *Proceedings of the 2021 International Conference on Management of Data*, pp. 577–590, 2021.

[26] Talia B Gillis, Vitaly Meursault, and Berk Ustun. Operationalizing the search for less discriminatory alternatives in fair lending. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 377–387, 2024.

[27] Michele E Gilman. Poverty lawgorithms: A poverty lawyer's guide to fighting automated decision-making harms on low-income communities. *Data & Society*, 2020.

[28] Sofie Goethals, David Martens, and Theodoros Evgeniou. Manipulation risks in explainable ai: The implications of the disagreement problem. *arXiv preprint arXiv:2306.13885*, 2023.

[29] Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. In *Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science*, pp. 111–122. ACM, 2016.

[30] Moritz Hardt, Eric Mazumdar, Celestine Mendler-Dünner, and Tijana Zrnic. Algorithmic collective action in machine learning. *arXiv preprint arXiv:2302.04262*, 2023.

[31] Mikella Hurley and Julius Adebayo. Credit scoring in the era of big data. *Yale JL & Tech.*, 18:148, 2016.

[32] Neil Jethani, Mukund Sudarshan, Ian Connick Covert, Su-In Lee, and Rajesh Ranganath. Fastshap: Real-time shapley value estimation. In *International conference on learning representations*, 2021.

[33] Hailey Joren, Chirag Nagpal, Katherine A Heller, and Berk Ustun. Participatory personalization in classification. *Advances in Neural Information Processing Systems*, 36:14113–14133, 2023.

[34] Hailey Joren, Charles Marx, and Berk Ustun. Classification with conceptual safeguards. *arXiv preprint arXiv:2411.04342*, 2024.

[35] Kaggle. Give Me Some Credit. http://www.kaggle.com/c/GiveMeSomeCredit/, 2011.

[36] Margot E Kaminski. The right to explanation, explained. In *Research Handbook on Information Law and Governance*, pp. 278–299. Edward Elgar Publishing, 2021.

[37] Amir-Hossein Karimi, Julius Von Kügelgen, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse under imperfect causal knowledge: a probabilistic approach. *Advances in neural information processing systems*, 33:265–277, 2020.

[38] Amir-Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse: from counterfactual explanations to interventions. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pp. 353–362, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445899. URL https://doi.org/10.1145/3442188.3445899.

[39] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. Interpreting interpretability: understanding data scientists' use of interpretability tools for machine learning. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pp. 1–14, 2020.

[40] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International Conference on Machine Learning*, pp. 5338–5348. PMLR, 2020.

[41] Seunghun Koh, Byung Hyung Kim, and Sungho Jo. Understanding the user perception and experience of interactive algorithmic recourse customization. *ACM Transactions on Computer-Human Interaction*, 2024.

[42] Avni Kothari, Bogdan Kulynych, Tsui-Wei Weng, and Berk Ustun. Prediction without preclusion: Recourse verification with reachable sets. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=SCQfYpdoGE.

[43] Satyapriya Krishna, Tessa Han, Alex Gu, Javin Pombra, Shahin Jabbari, Steven Wu, and Himabindu Lakkaraju. The disagreement problem in explainable machine learning: A practitioner's perspective. *arXiv preprint arXiv:2202.01602*, 2022.

[44] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *NeurIPS*, 2017.

[45] Charles Marx, Richard Phillips, Sorelle Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. Disentangling influence: Using disentangled representations to audit model predictions. *Advances in Neural Information Processing Systems*, 32, 2019.

[46] Charles Marx, Flavio Calmon, and Berk Ustun. Predictive multiplicity in classification. In *Proceedings of Machine Learning and Systems 2020*, pp. 9215–9224. 2020.

[47] Ramaravind K Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pp. 607–617, 2020.

[48] Bitya Neuhof and Yuval Benjamini. Confident feature ranking. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li (eds.), *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 1468–1476. PMLR, 02–04 May 2024. URL https://proceedings.mlr.press/v238/neuhof24a.html.

[49] Duy Nguyen, Ngoc Bui, and Viet Anh Nguyen. Distributionally robust recourse action. *arXiv preprint arXiv:2302.11211*, 2023.

[50] Martin Pawelczyk, Teresa Datta, Johan Van den Heuvel, Gjergji Kasneci, and Himabindu Lakkaraju. Probabilistically robust recourse: Navigating the trade-offs between costs and robustness in algorithmic recourse. In *The Eleventh International Conference on Learning Representations*, 2023.

[51] Manish Raghavan, Solon Barocas, Jon Kleinberg, and Karen Levy. Mitigating bias in algorithmic hiring: Evaluating claims and practices. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pp. 469–481, 2020.

[52] Marco Tulio Ribeiro, Sameer Singh, and Carlos 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*, pp. 1135–1144. ACM, 2016.

[53] Pedro Saleiro, Benedict Kuester, Loren Hinkson, Jesse London, Abby Stevens, Ari Anisfeld, Kit T Rodolfa, and Rayid Ghani. Aequitas: A bias and fairness audit toolkit. *arXiv preprint arXiv:1811.05577*, 2018.

[54] Andrew D Selbst and Solon Barocas. The intuitive appeal of explainable machines. 2018.

[55] Julian Skirzynski, David Danks, and Berk Ustun. Discrimination exposed? on the reliability of explanations for discrimination detection. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency*, 2025.

[56] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pp. 180–186, 2020.

[57] Winnie F Taylor. Meeting the equal credit opportunity act's specificity requirement: Judgmental and statistical scoring systems. *Buff. L. Rev.*, 29:73, 1980.

[58] The Lawyers' Committee for Civil Rights Under Law. Online civil rights act, December, 2023. URL https://www.lawyerscommittee.org/online-civil-rights-act.

[59] Sohini Upadhyay, Shalmali Joshi, and Himabindu Lakkaraju. Towards robust and reliable algorithmic recourse. *arXiv preprint arXiv:2102.13620*, 2021.

[60] Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT* '19, pp. 10–19. ACM, 2019. ISBN 978-1-4503-6125-5. doi: 10.1145/3287560.3287566.

[61] Suresh Venkatasubramanian and Mark Alfano. The philosophical basis of algorithmic recourse. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 284–293, 2020.

[62] Jamelle Watson-Daniels, David C. Parkes, and Berk Ustun. Predictive multiplicity in probabilistic classification. In *AAAI Conference on Artificial Intelligence*, 06 2023.

[63] White House. Blueprint for an AI bill of rights: Making automated systems work for the American people. The White House Office of Science and Technology Policy, October, 2022. URL https://www.whitehouse.gov/ostp/ai-bill-of-rights/.

[64] Laurence A Wolsey. *Integer programming*. John Wiley & Sons, 2020.

[65] S Wykstra. Government's use of algorithm serves up false fraud charges. undark, 6 january, 2020.

[66] Muhammad Rehman Zafar and Naimul Mefraz Khan. Dlime: A deterministic local interpretable model-agnostic explanations approach for computer-aided diagnosis systems. *arXiv preprint arXiv:1906.10263*, 2019.

[67] Zhengze Zhou, Giles Hooker, and Fei Wang. S-lime: Stabilized-lime for model explanation. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 2429–2438, 2021.

# Supplementary Materials

# A    SUPPLEMENTARY MATERIAL FOR SECTION 4

In this section, we provide additional implementation details for our methods to compute responsiveness scores in Section 4.

**Partitions**    Our implementation of algorithms in Section 4 partitions the feature space into disjoint sets. Each part is made up of features that share joint constraints. More formally, we define a partition $\{\pi_1, \pi_2, \ldots, \pi_k\}$ of $[d]$ such that given two parts $\pi_m, \pi_n$, there are no joint constraints between all pairs $(p, q) \in \pi_m \times \pi_n$ of features. Another way to think about feature partitions would be as connected components in a graph, where features are nodes and edges represent joint constraints (i.e., $\exists$ edge $(p, q) \iff$ there are joint actionability constraints between $p$ and $q$).

In what follows, we denote $\pi'$ as a part that contains $j$ (i.e., $j \in \pi'$).

## A.1    IMPLEMENTATION DETAILS FOR REACHABLE SET ENUMERATION

**Description of** Find1DAction **Routine**    The Find1DAction routine in Algorithm 2 enumerates a set of possible actions from an intervention on single feature by recovering all possible solutions for an optimization problem of the form:

$$\text{Find1DAction}(\boldsymbol{x}, A_j) := \operatorname*{argmin}_{\boldsymbol{a}} \|\boldsymbol{a}\|_1 \text{ s.t. } \boldsymbol{a} \in A_j(\boldsymbol{x})$$

The routine takes as input:

- $\boldsymbol{x} \in \mathcal{X}$, a point
- $A_j(\boldsymbol{x})$, the action set for feature $j$, representing both actionability constraints and feasible actions.
- $\mathcal{A}_j^{\text{opt}}$, a set of $[L] := |\mathcal{A}_j^{\text{opt}}|$ actions enumerated over previous iterations

Since actionability constraints specified in $A_j(\boldsymbol{x})$ precisely define feasible actions, we overload notation to allow $A_j(\boldsymbol{x})$ to represent the set of feasible actions.

At each iteration, it searches for the nearest single-feature action from the set $\boldsymbol{a} \in A_j(\boldsymbol{x})$ by solving a mixed-integer program formulation shown in Eq. (2). The optimization procedure returns the nearest action when it exists, or returns a certificate of infeasibility, which indicates that there are no more actions to enumerate. As detailed in Algorithm 2, if we find a solution $\boldsymbol{a}^*$, we remove it from $A_j(\boldsymbol{x})$ for the next iteration. In practice, we add solutions from each iteration to $\mathcal{A}_j^{\text{opt}}$ and add constraints with respect to each solution:

$$
\begin{aligned}
\min_{\boldsymbol{a}} \quad & \sum_{k \in \pi'} a_k^+ + a_k^- & & & \\
\text{s.t.} \quad & a_j \neq 0 & & \textit{intervene on } j & (2a) \\
& a_k = a_k^+ - a_k^- & k \in \pi' & \textit{reconstruction of } a_k & (2b) \\
& a_k^+ \geq a_k & k \in \pi' & \textit{positive component of } a_k & (2c) \\
& a_k^- \geq -a_k & k \in \pi' & \textit{negative component of } a_k & (2d) \\
& a_k^+ \leq \left| \sup_{\boldsymbol{a}' \in A_j(\boldsymbol{x})} a_k' \right| \sigma_k & k \in \pi' & a_k^+ > 0 \implies \sigma_k = 1 & (2e) \\
& a_k^- \leq \left| \inf_{\boldsymbol{a}' \in A_j(\boldsymbol{x})} a_k' \right| (1 - \sigma_k) & k \in \pi' & a_k^- > 0 \implies \sigma_k = 0 & (2f) \\
& a_k = a_{k,l} + \delta_{k,l}^+ - \delta_{k,l}^- & k \in \pi', \boldsymbol{a}_l \in \mathcal{A}_j^{\text{opt}} & \textit{maintain distance from prior actions} & (2g) \\
& \varepsilon_{\min} \leq \sum_{k \in \pi'} (\delta_{k,l}^+ + \delta_{k,l}^-) & \boldsymbol{a}_l \in \mathcal{A}_j^{\text{opt}} & \textit{any solution is } \varepsilon_{min} \textit{ away from } \boldsymbol{a}_l & (2h) \\
& \delta_{k,l}^+ \leq M_{k,l}^+ u_{k,l} & k \in \pi', l \in [L] & \delta_{k,l}^+ > 0 \implies u_{k,l} = 1 & (2i) \\
& \delta_{k,l}^- \leq M_{k,l}^- (1 - u_{k,l}) & k \in \pi', l \in [L] & \delta_{k,l}^- > 0 \implies u_{k,l} = 0 & (2j) \\
& \boldsymbol{a} \in A_j(\boldsymbol{x}) & & \textit{joint actionability constraints on } j & (2k) \\
& a_k^+, a_k^- \in \mathbb{R}_+ & k \in \pi' & \textit{absolute value of } a_k & (2l) \\
& \delta_{k,l}^+, \delta_{k,l}^- \in \mathbb{R}_+ & k \in \pi', l \in [L] & \textit{signed distances from } a_{k,l} & (2m) \\
& u_{k,l} \in \{0, 1\} & k \in \pi', l \in [L] & \textit{sign indicator of } \delta_{k,l} & (2n) \\
& \sigma_k \in \{0, 1\} & k \in \pi' & \textit{sign indicator of } a_k & (2o)
\end{aligned}
$$

This formulation finds the closest feasible action that at $\varepsilon_{\min}$ away from each action in the set $A_j(\boldsymbol{x})$. Here, the objective minimizes the $L_1$-norm of $a_k$ in terms of its positive and negative $a_k^+ - a_k^-$, which are defined in constraints (2b) and (2l). The constraints enforce a minimum distance between $a_k$ and the $l^{\text{th}}$ solution from the set $\mathcal{A}_j^{\text{opt}}$ in terms of the distance variables $\delta_{k,l}^+$ and $\delta_{k,l}^-$, which are defined in constraints (2g) and (2m). Here, $\sigma_k := \mathbb{1}[a_k > 0]$ and $u_{k,l} := \delta_{k,l}$ are binary variables set to 1 when $a_k$ and $\delta_{k,l}$ have positive signs, respectively. The formulation ensures these variables to ensure that signed components can have a positive value through constraints (2o) and (2n), respectively.

The first constraint, (2a) enforces that we intervene on feature $j$. The remaining constraints describe three key requirements for $\boldsymbol{a}$:

1. Sufficient distance from prior solutions (constraint (2h))

2. Adherence to separable actionability constraints (constraint (2e), (2f), (2i), (2j))

3. Adherence to joint actionability constraints (constraint (2k))

Constraint (2h) ensures that given $\varepsilon_{\min} > 0$, $\|\boldsymbol{a} - \boldsymbol{a}_l\|_1 \geq \varepsilon_{\min} \ \forall \ \boldsymbol{a}_l \in \mathcal{A}_j^{\text{opt}}$. We set $\varepsilon_{\min} = 0.5$ for our experiments with discrete datasets.

Constraints (2e), (2f) ensure that $a_k$ is feasible under separable constraints on $k \in \pi'$ and that only one of $a_k^+$ or $a_k^-$ is strictly positive. Similarly, constraints (2i), (2j) ensure that the distances between $\boldsymbol{a}$ and each $\boldsymbol{a}_l$ are within some bound. We achieve this by setting "Big-M" parameters $M_{k,l}^+, M_{k,l}^-$, which represent the upper bound for $\delta_{k,l}^+$ and $\delta_{k,l}^-$. For each feature $k \in \pi'$, we let

$$M_{k,l}^+ := \left| \sup_{\boldsymbol{a}' \in A_j(\boldsymbol{x})} a_k' - a_{k,l} \right|, \ M_{k,l}^- := \left| \inf_{\boldsymbol{a}' \in A_j(\boldsymbol{x})} a_k' - a_{k,l} \right|,$$

Along with the indicator variable $u_{k,l}$, $M_{k,l}^+, M_{k,l}^-$ ensure that only one of $\delta_{k,l}^+$ or $\delta_{k,l}^-$ is strictly positive and is feasible under separable actionability constraints.

Constraint (2k) ensures that $\boldsymbol{a}$ also adheres to joint actionability constraints. These constraints will exist if and only if $|\pi'| > 1$. See [42] for examples of how we can explicitly encode joint actionability constraints into Eq. (2).

This formulation is adapts the MIP in [42] for a task where we only need to enumerate actions with respect to a single-feature intervention $\boldsymbol{v}$.

## A.2 Implementation Details for Reachable Set Sampling

The sampling algorithm (Algorithm 1) requires additional considerations—most notably the sample size $N$.

**Choosing a Sample Size** The sample size $N$ controls the precision of the estimated responsiveness score $\hat{\mu}_j(\boldsymbol{x})$. We formalize precision using confidence intervals by treating $\hat{\mu}_j(\boldsymbol{x})$ as a binomial distribution parameter:

**Remark 2.** *Given a point $\boldsymbol{x} \in \mathcal{X}$, let $\hat{R}_j(\boldsymbol{x})$ denote a sample of $N$ points drawn uniformly from the reachable set $R_j(\boldsymbol{x})$. Given any model $f : \mathcal{X} \to \mathcal{Y}$, we can estimate the responsiveness score for feature $j$ as $\hat{\mu}_j(\boldsymbol{x}) := \frac{1}{N} \sum_{\boldsymbol{x}' \in \hat{R}_j(\boldsymbol{x})} \mathbb{1}[f(\boldsymbol{x}') = \hat{y}^{target}]$. Given a significance level $\alpha \in (0, 1)$, we have that:*

$$\Pr(\mu_j(\boldsymbol{x}) \in [\tilde{\mu}_j(\boldsymbol{x}) - \mathcal{E}, \tilde{\mu}_j(\boldsymbol{x}) + \mathcal{E}]) \geq 1 - \alpha$$

*Here: $\mathcal{E} := \kappa \sqrt{\frac{1}{N + \kappa^2} \tilde{\mu}_j(\boldsymbol{x})(1 - \tilde{\mu}_j(\boldsymbol{x}))}$ and $\tilde{\mu}_j(\boldsymbol{x}) := \frac{1}{N + \kappa^2}\left(S + \frac{\kappa^2}{2}\right)$ is a corrected estimator to improve coverage when $\mu_j(\boldsymbol{x}) \in \{0, 1\}$ [11], $S := |\{\boldsymbol{x}' \in \hat{R}_j(\boldsymbol{x}) \mid f(\boldsymbol{x}') = \hat{y}^{target}\}|$ is the subset of responsive points, and $\kappa := \Phi^{-1}(1 - \frac{\alpha}{2})$ is a constant based on the Normal CDF $\Phi(\cdot)$.*

The Agresti–Coull interval above is an approximate confidence interval for a binomial proportion [3], offering an improvement over the standard normal approximation known as the Wald Interval. It is particularly effective for small proportion values, providing more reliable coverage—the probability that the interval contains the true parameter value [11].

**Remark 3.** *Given $\alpha$ and $N$, $\mathcal{E}$ is maximized when $S = \frac{N}{2}$ and attains its minimum at $S = 0$ and $S = N$.*

*Proof.* Let $z = \frac{S + \frac{\kappa^2}{2}}{N + \kappa^2}$. Then, we have:

$$\mathcal{E} = \kappa \sqrt{\frac{z(1 - z)}{(N + \kappa^2)}}$$

18

Since $N$ and $\alpha$ are fixed ($\kappa$ is defined by $\alpha$), $\mathcal{E}$ can be expressed as a function of $z$ of the form $h(z) = c\sqrt{z(1-z)}$ where $c \in \mathbb{R}^+$. We observe that h(z) is a concave function whose first derivative can be expressed as:

$$h'(z) = \frac{c(1-2z)}{2\sqrt{z(1-z)}} \qquad h''(z) = -\frac{c}{4}[z(1-z)]^{-\frac{3}{2}}$$

Since $z > 0$ and $c > 0$, we can see that $h(z)$ attains a maximum value at $z' = \frac{1}{2}$ since $h(z') = 0$ and $h''(z) < 0$.

We see that $h'(z) > 0$ where $z < \frac{1}{2}$, meaning it is increasing for $z \in (0, \frac{1}{2}]$. Thus, the local minimum is achieved at the smallest possible $z$—when $S = 0$.

Similarly, for $z \in [\frac{1}{2}, 1)$, $h'(z) < 0$ and the local minimum is achieved at the largest possible $z$—when $S = N$.

Note that the value of $h$ (or $\mathcal{E}$) are the same at those two points.

$\square$

Using the Remark 3, we can a sample size $N$ in terms of $\alpha$ in following ways.

1. Control the precision when $S = 0$ (i.e., no points in $\hat{R}_j(\boldsymbol{x})$ are responsive) $\iff$ control the width of the shortest interval

2. Control the precision when $S = \frac{N}{2}$ (i.e., half of the points in $\hat{R}_j(\boldsymbol{x})$ are responsive) $\iff$ control the width of the widest interval

Either way, we fix $\alpha$ and solve for $N$ given the width of the interval $\mathcal{E}$ at a specified $S$. Below we provide a table of the smallest $N$ needed for different $\mathcal{E}$—interval widths—at common values of $\alpha$ for the two methods:

| | Width of Interval ($\mathcal{E}$) | | | |
|---|---|---|---|---|
| $\alpha$ | 0.01 | 0.02 | 0.05 | 0.10 |
| 0.01 | 461 | 227 | 86 | 39 |
| 0.05 | 267 | 132 | 50 | 23 |
| 0.10 | 188 | 93 | 35 | 16 |

Table 6: Minimum $N$ required to ensure the shortest confidence interval is less than $2\mathcal{E}$ (Method 1)

| | Width of Interval ($\mathcal{E}$) | | | |
|---|---|---|---|---|
| $\alpha$ | 0.01 | 0.02 | 0.05 | 0.10 |
| 0.01 | 16581 | 4141 | 657 | 160 |
| 0.05 | 9600 | 2398 | 381 | 93 |
| 0.10 | 6762 | 1689 | 268 | 65 |

Table 7: Minimum $N$ required to ensure the widest confidence interval is less than $2\mathcal{E}$ (Method 2)

**Description of the Sample1DAction Routine** Let $j$ be the feature that we are intervening on.

**Case 1:** $|\pi'| = 1$ (i.e., $\pi' = \{j\}$, $j$ is not jointly constrained with other features).

Here, there are no downstream effects from intervening on $j$. We take a uniformly random intervention from $V_j(\boldsymbol{x})$:

$$\boldsymbol{a}^* \sim A_j(\boldsymbol{x})$$

which abides by $j$'s separable actionability constraints like feature bounds and monotonicity.

**Case 2:** $|\pi'| > 1$ (i.e., $j$ is jointly constrained with other features)

We breakdown the partition $\pi'$ into three disjoint subsets:

$$\pi' = \{j\} \cup \pi'_{\text{disc}} \cup \pi'_{\text{cts}}$$

where $\pi'_{\text{disc}}$ and $\pi'_{\text{cts}}$ are the sets of discrete and and continuous features in $\pi'$ respectively.

We consider the following three sub-cases:

*Case 2a:* $|\pi'_{\text{cts}}| = 0$—all features in $\pi'$ are discrete.

We repeatedly solve the MIP in Find1DAction and take a sample from the resulting set of feasible actions.

*Case 2b:* $|\pi'_{\text{disc}}| = 0$—all features in $\pi'$ are continuous.

We sample action values that abide by separable actionability constraints for each feature in $\pi'_{\text{disc}}$.

*Case 2c:* $|\pi'_{\text{cts}}|, |\pi'_{\text{disc}}| > 0$—part contains discrete and continuous features.

19

We run the sampling steps in *Case 2a, 2b* for $\pi'_{\text{disc}} \cup \{j\}$ and $\pi'_{\text{cts}}$ to get $\boldsymbol{a}_{\text{disc}}$ and $\boldsymbol{a}_{\text{cts}}$.

We then check feasibility on $\boldsymbol{a}^* = \boldsymbol{a}_{\text{disc}} + \boldsymbol{a}_{\text{cts}}$ by running $\mathsf{CheckFeasibility}(\boldsymbol{a}^*, A_j(\boldsymbol{x}))$.

**Description of $\mathsf{CheckFeasibility}$ Routine** We describe the implementation for the $\mathsf{CheckFeasibility}(\boldsymbol{x}, \boldsymbol{a}^*, A_j)$ in Algorithm 1. Contrary to the MIP formulation in Section A.1, given the original point $\boldsymbol{x} \in \mathcal{X}$ and the sampled action $\boldsymbol{a}^*$, we solve the MIP once.

$$\min_{\boldsymbol{a}} \quad 1$$

$$
\begin{aligned}
\text{s.t.} \quad & \boldsymbol{a} = \boldsymbol{a}^* && && \textit{match action } \boldsymbol{a}^* && \text{(3a)}\\
& a_k = a_k^+ - a_k^- && k \in \pi' && \textit{reconstruction of } a_k && \text{(3b)}\\
& a_k^+ \geq a_k && k \in \pi' && \textit{positive component of } a_k && \text{(3c)}\\
& a_k^- \geq -a_k && k \in \pi' && \textit{negative component of } a_k && \text{(3d)}\\
& a_k^+ \leq \left| \sup_{\boldsymbol{a}' \in A_j(\boldsymbol{x})} a_k' \right| \sigma_k && k \in \pi' && a_k^+ > 0 \implies \sigma_k = 1 && \text{(3e)}\\
& a_k^- \leq \left| \inf_{\boldsymbol{a}' \in A_j(\boldsymbol{x})} a_k' \right| (1 - \sigma_k) && k \in \pi' && a_k^- > 0 \implies \sigma_k = 0 && \text{(3f)}\\
& \boldsymbol{a} \in A_j(\boldsymbol{x}) && && \textit{joint actionability constraints on } j && \text{(3g)}\\
& a_k^+, a_k^- \in \mathbb{R}_+ && k \in \pi' && \textit{absolute value of } a_k && \text{(3h)}\\
& \sigma_k \in \{0, 1\} && k \in \pi' && \textit{sign indicator of } a_k && \text{(3i)}
\end{aligned}
$$

The formulation is a variant of the problem in Section A.1, where:

- $\boldsymbol{a} = \boldsymbol{a}^*$,
- $\mathcal{A}_j^{\text{opt}} = \varnothing$,
- and set the objective to $\min_{\boldsymbol{a}} \; 1$

Hence $\mathsf{CheckFeasibility}(\boldsymbol{x}, \boldsymbol{a}^*, A_j) = 1$ if $\boldsymbol{a}^*$ is feasible under actionability constraints and 0 otherwise.

In practice, we run $\mathsf{Sample1DAction}$ and $\mathsf{CheckFeasibility}$ in Algorithm 1 in batches for efficiency; rather than sampling one action and checking feasibility, we sample $\tilde{N} >> N$ points and then check feasibility at once. We sample more than the required $N$ points to account for rejected samples in the $\mathsf{CheckFeasibility}$ step.

## B  SUPPLEMENTARY EXPERIMENT DETAILS

### B.1  DETAILS FOR THE `heloc` DATASET

The `heloc` dataset was created to predict repayment on Home Equity Line of Credit HELOC applications; these are loans that use people's homes as collateral. The anonymized version of the dataset was developed by FICO for use in an Explainable Machine Learning Challenge in 2018 [22]. Each instance in the dataset is an application for a home equity loan from a single applicant. There are $n = 10,459$ instances and $d = 23$ features. Here, the label $y_i = 0$ if an applicant has been more than 90 days overdue on their payments in the last 2 years and $y_i = 1$ otherwise. We thermometer encode continuous or integer features after dropping rows and features with missing data (see Table 8). See GitHub for dataset processing code.

| Name | Type | LB | UB | Actionability | Sign | Joint Constraints | Partition ID |
|---|---|---|---|---|---|---|---|
| NumInstallTrades≥2 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTradesWBalance≥2 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTrades≥3 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTradesWBalance≥3 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTrades≥5 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTradesWBalance≥5 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTrades≥7 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumInstallTradesWBalance≥7 | $\{0,1\}$ | 0 | 1 | Yes | + | 20, 21, 24, 25, 28, 29, 32, 33 | 14 |
| NumRevolvingTrades≥2 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTradesWBalance≥2 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTrades≥3 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTradesWBalance≥3 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTrades≥5 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTradesWBalance≥5 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTrades≥7 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| NumRevolvingTradesWBalance≥7 | $\{0,1\}$ | 0 | 1 | Yes | + | 22, 23, 26, 27, 30, 31, 34, 35 | 15 |
| YearsOfAccountHistory | $\mathbb{Z}$ | 0 | 50 | No | | 5, 17, 18, 19 | 5 |
| YearsSinceLastDelqTrade≤1 | $\{0,1\}$ | 0 | 1 | Yes | + | 5, 17, 18, 19 | 5 |
| YearsSinceLastDelqTrade≤3 | $\{0,1\}$ | 0 | 1 | Yes | + | 5, 17, 18, 19 | 5 |
| YearsSinceLastDelqTrade≤5 | $\{0,1\}$ | 0 | 1 | Yes | + | 5, 17, 18, 19 | 5 |
| NetFractionInstallBurden≥10 | $\{0,1\}$ | 0 | 1 | Yes | + | 36, 37, 38 | 16 |
| NetFractionInstallBurden≥20 | $\{0,1\}$ | 0 | 1 | Yes | + | 36, 37, 38 | 16 |
| NetFractionInstallBurden≥50 | $\{0,1\}$ | 0 | 1 | Yes | + | 36, 37, 38 | 16 |
| NetFractionRevolvingBurden≥10 | $\{0,1\}$ | 0 | 1 | Yes | + | 39, 40, 41 | 17 |
| NetFractionRevolvingBurden≥20 | $\{0,1\}$ | 0 | 1 | Yes | + | 39, 40, 41 | 17 |
| NetFractionRevolvingBurden≥50 | $\{0,1\}$ | 0 | 1 | Yes | + | 39, 40, 41 | 17 |
| AvgYearsInFile≥3 | $\{0,1\}$ | 0 | 1 | Yes | + | 6, 7, 8 | 6 |
| AvgYearsInFile≥5 | $\{0,1\}$ | 0 | 1 | Yes | + | 6, 7, 8 | 6 |
| AvgYearsInFile≥7 | $\{0,1\}$ | 0 | 1 | Yes | + | 6, 7, 8 | 6 |
| MostRecentTradeWithinLastYear | $\{0,1\}$ | 0 | 1 | Yes | | 9, 10 | 7 |
| MostRecentTradeWithinLast2Years | $\{0,1\}$ | 0 | 1 | Yes | | 9, 10 | 7 |
| ExternalRiskEstimate≥40 | $\{0,1\}$ | 0 | 1 | No | | – | 0 |
| ExternalRiskEstimate≥50 | $\{0,1\}$ | 0 | 1 | No | | – | 1 |
| ExternalRiskEstimate≥60 | $\{0,1\}$ | 0 | 1 | No | | – | 2 |
| ExternalRiskEstimate≥70 | $\{0,1\}$ | 0 | 1 | No | | – | 3 |
| ExternalRiskEstimate≥80 | $\{0,1\}$ | 0 | 1 | No | | – | 4 |
| AnyDerogatoryComment | $\{0,1\}$ | 0 | 1 | No | | – | 8 |
| AnyTrade120DaysDelq | $\{0,1\}$ | 0 | 1 | No | | – | 9 |
| AnyTrade90DaysDelq | $\{0,1\}$ | 0 | 1 | No | | – | 10 |
| AnyTrade60DaysDelq | $\{0,1\}$ | 0 | 1 | No | | – | 11 |
| AnyTrade30DaysDelq | $\{0,1\}$ | 0 | 1 | No | | – | 12 |
| NoDelqEver | $\{0,1\}$ | 0 | 1 | No | | – | 13 |
| NumBank2NatlTradesWHighUtilizationGeq2 | $\{0,1\}$ | 0 | 1 | Yes | + | – | 18 |

**Table 8:** Separable Actionability Constraints for the processed `heloc` dataset. **Type** indicates the feature type ($\mathbb{Z}$ for integer, $\{0,1\}$ for binary). **LB**, **UB** are the lower and upper bounds for the feature. **Actionability** indicates whether the feature is globally actionable. **Sign** indicates if the feature can only increase (+) or decrease (-). **Joint Constraints** are a list non-separable constraint indices it is tied to (if any). **Partition ID** indicates which partition the feature belongs to.

**Actionability Constraints**  The joint actionability constraints listed in Table 8 include:

1. DirectionalLinkage: Actions on `NumRevolvingTradesWBalance≥2` will induce to actions on ['NumRevolvingTrades≥2'].Each unit change in `NumRevolvingTradesWBalance≥2` leads to:1.00-unit change in `NumRevolvingTrades≥2`

2. DirectionalLinkage: Actions on `NumInstallTradesWBalance≥2` will induce to actions on ['NumInstallTrades≥2'].Each unit change in `NumInstallTradesWBalance≥2` leads to:1.00-unit change in `NumInstallTrades≥2`

3. DirectionalLinkage: Actions on `NumRevolvingTradesWBalance≥3` will induce to actions on ['NumRevolvingTrades≥3'].Each unit change in `NumRevolvingTradesWBalance≥3` leads to:1.00-unit change in `NumRevolvingTrades≥3`

4. DirectionalLinkage: Actions on `NumInstallTradesWBalance≥3` will induce to actions on ['NumInstallTrades≥3'].Each unit change in `NumInstallTradesWBalance≥3` leads to:1.00-unit change in `NumInstallTrades≥3`

5. DirectionalLinkage: Actions on `NumRevolvingTradesWBalance`$\geq$5 will induce to actions on ['NumRevolvingTrades$\geq$5'].Each unit change in `NumRevolvingTradesWBalance`$\geq$5 leads to:1.00-unit change in `NumRevolvingTrades`$\geq$5

6. DirectionalLinkage: Actions on `NumInstallTradesWBalance`$\geq$5 will induce to actions on ['NumInstallTrades$\geq$5'].Each unit change in `NumInstallTradesWBalance`$\geq$5 leads to:1.00-unit change in `NumInstallTrades`$\geq$5

7. DirectionalLinkage: Actions on `NumRevolvingTradesWBalance`$\geq$7 will induce to actions on ['NumRevolvingTrades$\geq$7'].Each unit change in `NumRevolvingTradesWBalance`$\geq$7 leads to:1.00-unit change in `NumRevolvingTrades`$\geq$7

8. DirectionalLinkage: Actions on `NumInstallTradesWBalance`$\geq$7 will induce to actions on ['NumInstallTrades$\geq$7'].Each unit change in `NumInstallTradesWBalance`$\geq$7 leads to:1.00-unit change in `NumInstallTrades`$\geq$7

9. DirectionalLinkage: Actions on `YearsSinceLastDelqTrade`$\leq$1 will induce to actions on ['YearsOfAccountHistory'].Each unit change in `YearsSinceLastDelqTrade`$\leq$1 leads to:-1.00-unit change in `YearsOfAccountHistory`

10. DirectionalLinkage: Actions on `YearsSinceLastDelqTrade`$\leq$3 will induce to actions on ['YearsOfAccountHistory'].Each unit change in `YearsSinceLastDelqTrade`$\leq$3 leads to:-3.00-unit change in `YearsOfAccountHistory`

11. DirectionalLinkage: Actions on `YearsSinceLastDelqTrade`$\leq$5 will induce to actions on ['YearsOfAccountHistory'].Each unit change in `YearsSinceLastDelqTrade`$\leq$5 leads to:-5.00-unit change in `YearsOfAccountHistory`

12. ReachabilityConstraint: The values of [`MostRecentTradeWithinLastYear`, `MostRecentTradeWithinLast2Years`] must belong to one of 4 values with custom reachability conditions.

13. ThermometerEncoding: Actions on [`YearsSinceLastDelqTrade`$\leq$1, `YearsSinceLastDelqTrade`$\leq$3, `YearsSinceLastDelqTrade`$\leq$5] must preserve thermometer encoding of YearsSinceLastDelqTradeleq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `YearsSinceLastDelqTrade`$\leq$1 is the lowest-level dummy and `YearsSinceLastDelqTrade`$\leq$5 is the highest-level-dummy.

14. ThermometerEncoding: Actions on [`AvgYearsInFile`$\geq$3, `AvgYearsInFile`$\geq$5, `AvgYearsInFile`$\geq$7] must preserve thermometer encoding of AvgYearsInFilegeq., which can only increase.Actions can only turn on higher-level dummies that are off, where `AvgYearsInFile`$\geq$3 is the lowest-level dummy and `AvgYearsInFile`$\geq$7 is the highest-level-dummy.

15. ThermometerEncoding: Actions on [`NetFractionRevolvingBurden`$\geq$10, `NetFractionRevolvingBurden`$\geq$20, `NetFractionRevolvingBurden`$\geq$50] must preserve thermometer encoding of NetFractionRevolvingBurdengeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `NetFractionRevolvingBurden`$\geq$10 is the lowest-level dummy and `NetFractionRevolvingBurden`$\geq$50 is the highest-level-dummy.

16. ThermometerEncoding: Actions on [`NetFractionInstallBurden`$\geq$10, `NetFractionInstallBurden`$\geq$20, `NetFractionInstallBurden`$\geq$50] must preserve thermometer encoding of NetFractionInstallBurdengeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `NetFractionInstallBurden`$\geq$10 is the lowest-level dummy and `NetFractionInstallBurden`$\geq$50 is the highest-level-dummy.

17. ThermometerEncoding: Actions on [`NumRevolvingTradesWBalance`$\geq$2, `NumRevolvingTradesWBalance`$\geq$3, `NumRevolvingTradesWBalance`$\geq$5, `NumRevolvingTradesWBalance`$\geq$7] must preserve thermometer encoding of NumRevolvingTradesWBalancegeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `NumRevolvingTradesWBalance`$\geq$2 is the lowest-level dummy and `NumRevolvingTradesWBalance`$\geq$7 is the highest-level-dummy.

18. ThermometerEncoding: Actions on [`NumRevolvingTrades`$\geq$2, `NumRevolvingTrades`$\geq$3, `NumRevolvingTrades`$\geq$5, `NumRevolvingTrades`$\geq$7] must preserve thermometer encoding of NumRevolvingTradesgeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `NumRevolvingTrades`$\geq$2 is the lowest-level dummy and `NumRevolvingTrades`$\geq$7 is the highest-level-dummy.

19. ThermometerEncoding: Actions on [`NumInstallTradesWBalance`$\geq$2, `NumInstallTradesWBalance`$\geq$3, `NumInstallTradesWBalance`$\geq$5, `NumInstallTradesWBalance`$\geq$7] must preserve thermometer encoding of NumInstallTradesWBalancegeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `NumInstallTradesWBalance`$\geq$2 is the lowest-level dummy and `NumInstallTradesWBalance`$\geq$7 is the highest-level-dummy.

20. ThermometerEncoding: Actions on [`NumInstallTrades`$\geq$`2`, `NumInstallTrades`$\geq$`3`, `NumInstallTrades`$\geq$`5`, `NumInstallTrades`$\geq$`7`] must preserve thermometer encoding of NumInstallTradesgeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where `NumInstallTrades`$\geq$`2` is the lowest-level dummy and `NumInstallTrades`$\geq$`7` is the highest-level-dummy.

## B.2 DETAILS FOR THE german DATASET

The german dataset was created in 1994 and contains information about loan history, demographics, occupation, payment history, and whether or not somebody is a good customer [16]. Each instance is credit applicant. There are $n = 1,000$ instances and $d = 20$ features. The features are all either categorical or discrete. The label a indicates is a loan applicant is "good" ($y_i = 1$) or "bad" ($y_i = 0$). There are no missing values in the dataset. We renamed some of the features to be indicative of the values they represent. The dataset is self-contained and anonymous, and it includes features describing gender, age, and marital status.

| Name | Type | LB | UB | Actionability | Sign | Joint Constraints | Partition ID |
|---|---|---|---|---|---|---|---|
| Age | $\mathbb{Z}$ | 19 | 75 | No | | 0, 4, 12 | 0 |
| YearsAtResidence | $\mathbb{Z}$ | 0 | 7 | Yes | + | 0, 4, 12 | 0 |
| YearsEmployed≥1 | {0, 1} | 0 | 1 | Yes | + | 0, 4, 12 | 0 |
| CheckingAcct_exists | {0, 1} | 0 | 1 | Yes | + | 32, 33 | 30 |
| CheckingAcct≥0 | {0, 1} | 0 | 1 | Yes | + | 32, 33 | 30 |
| SavingsAcct_exists | {0, 1} | 0 | 1 | Yes | + | 34, 35 | 31 |
| SavingsAcct≥100 | {0, 1} | 0 | 1 | Yes | + | 34, 35 | 31 |
| Male | {0, 1} | 0 | 1 | No | | – | 1 |
| Single | {0, 1} | 0 | 1 | No | | – | 2 |
| ForeignWorker | {0, 1} | 0 | 1 | No | | – | 3 |
| LiablePersons | $\mathbb{Z}$ | 1 | 2 | No | | – | 4 |
| Housing=Renter | {0, 1} | 0 | 1 | No | | – | 5 |
| Housing=Owner | {0, 1} | 0 | 1 | No | | – | 6 |
| Housing=Free | {0, 1} | 0 | 1 | No | | – | 7 |
| Job=Unskilled | {0, 1} | 0 | 1 | No | | – | 8 |
| Job=Skilled | {0, 1} | 0 | 1 | No | | – | 9 |
| Job=Management | {0, 1} | 0 | 1 | No | | – | 10 |
| CreditAmt≥1000K | {0, 1} | 0 | 1 | No | | – | 11 |
| CreditAmt≥2000K | {0, 1} | 0 | 1 | No | | – | 12 |
| CreditAmt≥5000K | {0, 1} | 0 | 1 | No | | – | 13 |
| CreditAmt≥10000K | {0, 1} | 0 | 1 | No | | – | 14 |
| LoanDuration≤6 | {0, 1} | 0 | 1 | No | | – | 15 |
| LoanDuration≥12 | {0, 1} | 0 | 1 | No | | – | 16 |
| LoanDuration≥24 | {0, 1} | 0 | 1 | No | | – | 17 |
| LoanDuration≥36 | {0, 1} | 0 | 1 | No | | – | 18 |
| LoanRate | $\mathbb{Z}$ | 1 | 4 | No | | – | 19 |
| HasGuarantor | {0, 1} | 0 | 1 | Yes | + | – | 20 |
| LoanRequiredForBusiness | {0, 1} | 0 | 1 | No | | – | 21 |
| LoanRequiredForEducation | {0, 1} | 0 | 1 | No | | – | 22 |
| LoanRequiredForCar | {0, 1} | 0 | 1 | No | | – | 23 |
| LoanRequiredForHome | {0, 1} | 0 | 1 | No | | – | 24 |
| NoCreditHistory | {0, 1} | 0 | 1 | No | | – | 25 |
| HistoryOfLatePayments | {0, 1} | 0 | 1 | No | | – | 26 |
| HistoryOfDelinquency | {0, 1} | 0 | 1 | No | | – | 27 |
| HistoryOfBankInstallments | {0, 1} | 0 | 1 | Yes | + | – | 28 |
| HistoryOfStoreInstallments | {0, 1} | 0 | 1 | Yes | + | – | 29 |

**Table 9:** Separable Actionability Constraints for the processed german dataset. **Type** indicates the feature type ($\mathbb{Z}$ for integer, {0, 1} for binary). **LB**, **UB** are the lower and upper bounds for the feature. **Actionability** indicates whether the feature is globally actionable. **Sign** indicates if the feature can only increase (+) or decrease (-). **Joint Constraints** are a list non-separable constraint indices it is tied to (if any). **Partition ID** indicates which partition the feature belongs to.

**Actionability Constraints** The joint actionability constraints listed in Table 9 include:

1. DirectionalLinkage: Actions on YearsAtResidence will induce to actions on ['Age'].Each unit change in YearsAtResidence leads to:1.00-unit change in Age

2. DirectionalLinkage: Actions on YearsEmployed≥1 will induce to actions on ['Age'].Each unit change in YearsEmployed≥1 leads to:1.00-unit change in Age

3. ThermometerEncoding: Actions on [CheckingAcctexists, CheckingAcct≥0] must preserve thermometer encoding of CheckingAcct., which can only increase.Actions can only turn on higher-level dummies that are off, where CheckingAcctexists is the lowest-level dummy and CheckingAcct≥0 is the highest-level-dummy.

4. ThermometerEncoding: Actions on [SavingsAcctexists, SavingsAcct≥100] must preserve thermometer encoding of SavingsAcct., which can only increase.Actions can only turn on higher-level dummies that are off, where SavingsAcctexists is the lowest-level dummy and SavingsAcct≥100 is the highest-level-dummy.

## B.3    DETAILS FOR THE `givemecredit` DATASET

The `givemecredit` dataset is used to determine whether a loan should be given or denied [35]. The label indicates whether someone was 90 days past due in the two years following data collection. Delinquency refers to a debt with an overdue payment; this dataset is used to predict if someone will experience financial distress in the next two years.It contains information about $n = 120,268$ loan recipients, and each instance represents a borrower. There are $d = 10$ features before preprocessing. Here, the label is $y_i = 0$ if an applicant has had a serious delinquency in two years and $y_i$ otherwise. The data is self-contained and anonymous, and it contains features describing age, income, and the number of dependents.

| Name | Type | LB | UB | Actionability | Sign | Joint Constraints | Partition ID |
|---|---|---|---|---|---|---|---|
| CreditLineUtilization≥10.0 | $\{0,1\}$ | 0 | 1 | Yes | | 12, 13, 14, 15, 16 | 10 |
| CreditLineUtilization≥20.0 | $\{0,1\}$ | 0 | 1 | Yes | | 12, 13, 14, 15, 16 | 10 |
| CreditLineUtilization≥50.0 | $\{0,1\}$ | 0 | 1 | Yes | | 12, 13, 14, 15, 16 | 10 |
| CreditLineUtilization≥70.0 | $\{0,1\}$ | 0 | 1 | Yes | | 12, 13, 14, 15, 16 | 10 |
| CreditLineUtilization≥100.0 | $\{0,1\}$ | 0 | 1 | Yes | | 12, 13, 14, 15, 16 | 10 |
| MonthlyIncome≥3K | $\{0,1\}$ | 0 | 1 | Yes | + | 9, 10, 11 | 9 |
| MonthlyIncome≥5K | $\{0,1\}$ | 0 | 1 | Yes | + | 9, 10, 11 | 9 |
| MonthlyIncome≥10K | $\{0,1\}$ | 0 | 1 | Yes | + | 9, 10, 11 | 9 |
| AnyRealEstateLoans | $\{0,1\}$ | 0 | 1 | Yes | + | 17, 18 | 11 |
| MultipleRealEstateLoans | $\{0,1\}$ | 0 | 1 | Yes | + | 17, 18 | 11 |
| AnyCreditLinesAndLoans | $\{0,1\}$ | 0 | 1 | Yes | + | 19, 20 | 12 |
| MultipleCreditLinesAndLoans | $\{0,1\}$ | 0 | 1 | Yes | + | 19, 20 | 12 |
| Age≤24 | $\{0,1\}$ | 0 | 1 | No | − | | 0 |
| Age_bt_25_to_30 | $\{0,1\}$ | 0 | 1 | No | − | | 1 |
| Age_bt_30_to_59 | $\{0,1\}$ | 0 | 1 | No | − | | 2 |
| Age≥60 | $\{0,1\}$ | 0 | 1 | No | − | | 3 |
| NumberOfDependents=0 | $\{0,1\}$ | 0 | 1 | No | − | | 4 |
| NumberOfDependents=1 | $\{0,1\}$ | 0 | 1 | No | − | | 5 |
| NumberOfDependents≥2 | $\{0,1\}$ | 0 | 1 | No | − | | 6 |
| NumberOfDependents≥5 | $\{0,1\}$ | 0 | 1 | No | − | | 7 |
| DebtRatio≥1 | $\{0,1\}$ | 0 | 1 | Yes | + | − | 8 |
| HistoryOfLatePayment | $\{0,1\}$ | 0 | 1 | No | − | | 13 |
| HistoryOfDelinquency | $\{0,1\}$ | 0 | 1 | No | − | | 14 |

**Table 10:** Separable Actionability Constraints for the processed `givemecredit` dataset. **Type** indicates the feature type ($\mathbb{Z}$ for integer, $\{0,1\}$ for binary). **LB**, **UB** are the lower and upper bounds for the feature. **Actionability** indicates whether the feature is globally actionable. **Sign** indicates if the feature can only increase (+) or decrease (-). **Joint Constraints** are a list non-separable constraint indices it is tied to (if any). **Partition ID** indicates which partition the feature belongs to.

**Actionability Constraints**    The joint actionability constraints listed in Table 10 include:

1. ThermometerEncoding:    Actions    on    [MonthlyIncome≥3K,    MonthlyIncome≥5K, MonthlyIncome≥10K] must preserve thermometer encoding of MonthlyIncomegeq., which can only increase.Actions can only turn on higher-level dummies that are off, where MonthlyIncome≥3K is the lowest-level dummy and MonthlyIncome≥10K is the highest-level-dummy.

2. ThermometerEncoding:    Actions    on    [CreditLineUtilization≥10.0, CreditLineUtilization≥20.0,    CreditLineUtilization≥50.0, CreditLineUtilization≥70.0, CreditLineUtilization≥100.0] must preserve thermometer encoding of CreditLineUtilizationgeq., which can only decrease.Actions can only turn off higher-level dummies that are on, where CreditLineUtilization≥10.0 is the lowest-level dummy and CreditLineUtilization≥100.0 is the highest-level-dummy.

3. ThermometerEncoding: Actions on [AnyRealEstateLoans, MultipleRealEstateLoans] must preserve thermometer encoding of continuousattribute., which can only decrease.Actions can only turn off higher-level dummies that are on, where AnyRealEstateLoans is the lowest-level dummy and MultipleRealEstateLoans is the highest-level-dummy.

4. ThermometerEncoding:    Actions    on    [AnyCreditLinesAndLoans, MultipleCreditLinesAndLoans] must preserve thermometer encoding of continuousattribute., which can only decrease.Actions can only turn off higher-level dummies that are on, where AnyCreditLinesAndLoans is the lowest-level dummy and MultipleCreditLinesAndLoans is the highest-level-dummy.

## B.4 OVERVIEW OF MODEL PERFORMANCE

We include the performance of the classifiers used in Section 5.

| | LR | | XGB | | RF | |
|---|---|---|---|---|---|---|
| Dataset | Train | Test | Train | Test | Train | Test |
| `heloc` $n = 5,842$ $d = 43$ FICO [22] | 0.772 | 0.788 | 0.859 | 0.785 | 0.780 | 0.790 |
| `german` $n = 1,000$ $d = 36$ Dua & Graff [16] | 0.819 | 0.760 | 0.971 | 0.794 | 0.828 | 0.766 |
| `givemecredit` $n = 120,268$ $d = 23$ Kaggle [35] | 0.841 | 0.844 | 0.875 | 0.793 | 0.864 | 0.835 |

**Table 11:** Train and Test AUC for models across all datasets. We optimized the model's hyperparameters through randomized search and divided the data into training and testing sets at an 80% and 20% ratio.

## C SUPPLEMENTARY EXPERIMENT RESULTS

### C.1 RESPONSIVENESS OF EXPLANATIONS FOR RANDOM FORESTS

| | | All Features | | Actionable Features | | |
|---|---|---|---|---|---|---|
| Dataset | Metrics | LIME | SHAP | LIME-AW | SHAP-AW | RESP |
| `heloc` $n = 5,842$ $d = 43$ features $d_A = 31$ mutable FICO [22] | % Presented with Explanations ↳ % All Features Unresponsive ↳ % At Least 1 Feature Responsive ↳ % All Features Responsive ↳ # Features Highlighted | 100.0% 86.5% 13.5% 0.0% 4.0 | 100.0% 78.2% 21.8% 0.0% 4.0 | 100.0% 77.1% 22.9% 0.0% 4.0 | 100.0% 76.7% 23.3% 0.5% 4.0 | 31.7% 0.0% 100.0% **100.0%** 2.4 |
| `german` $n = 1,000$ $d = 36$ features $d_A = 9$ mutable Dua & Graff [16] | % Presented with Explanations ↳ % All Features Unresponsive ↳ % At Least 1 Feature Responsive ↳ % All Features Responsive ↳ # Features Highlighted | 100.0% 100.0% 0.0% 0.0% 4.0 | 100.0% 89.1% 10.9% 0.0% 4.0 | 100.0% 76.6% 23.4% 0.0% 4.0 | 100.0% 64.6% 35.4% 0.0% 4.0 | 48.0% 0.0% 100.0% **100.0%** 2.2 |
| `givemecredit` $n = 120,268$ $d = 23$ features $d_A = 13$ mutable Kaggle [35] | % Presented with Explanations ↳ % All Features Unresponsive ↳ % At Least 1 Feature Responsive ↳ % All Features Responsive ↳ # Features Highlighted | 100.0% 56.5% 43.5% 0.0% 4.0 | 100.0% 26.8% 73.2% 0.5% 4.0 | 100.0% 28.4% 71.6% 1.4% 4.0 | 100.0% 21.0% 79.0% 11.4% 4.0 | 93.2% 0.0% 100.0% **100.0%** 2.9 |

**Table 12:** Responsiveness of feature-highlighting explanations for RF for all methods and datasets. We generate explanations that highlight up to 4 top-scoring features from a given method. We report the proportion of individuals receiving an explanation (*% Presented with Explanations*) and the mean number of features in each explanation (*# Features Highlighted*). We also show the proportion of instances where all features are unresponsive (*% All Features Unresponsive*) highlighting positive values; at least one feature is responsive (*% At Least 1 Feature Responsive*), or all features are responsive (*% All Features Responsive*) highlighting the **best value**.

## C.2 FEATURE RESPONSIVENESS RANKINGS

We include a plot to show how responsive features are at different rankings by LIME, SHAP, LIME-AW, SHAP-AW and RESP for each dataset. For every denied individual, we rank features by their absolute feature importance score returned by these methods. We exclude features with 0 attribution from the rankings.

The plots below show the % of times where the feature at each rank are responsive (i.e., feature has RESP > 0). It allows us to visualize and compare how often these methods assign high attribution to responsive features.
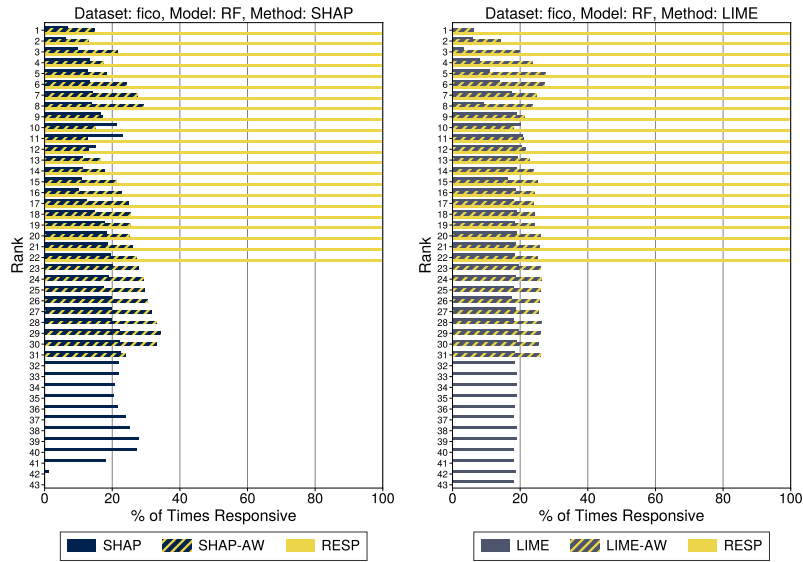
### C.2.1 `heloc`



**Figure 4:** Responsiveness of features for individuals who are denied credit by the LR model on the `fico` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."

**Figure 5:** Responsiveness of features for individuals who are denied credit by the XGB model on the `fico` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."
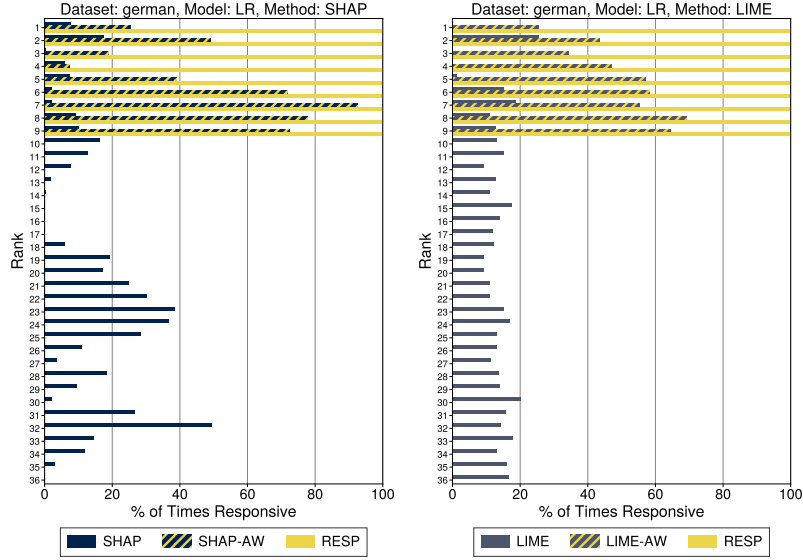


**Figure 6:** Responsiveness of features for individuals who are denied credit by the RF model on the `fico` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."
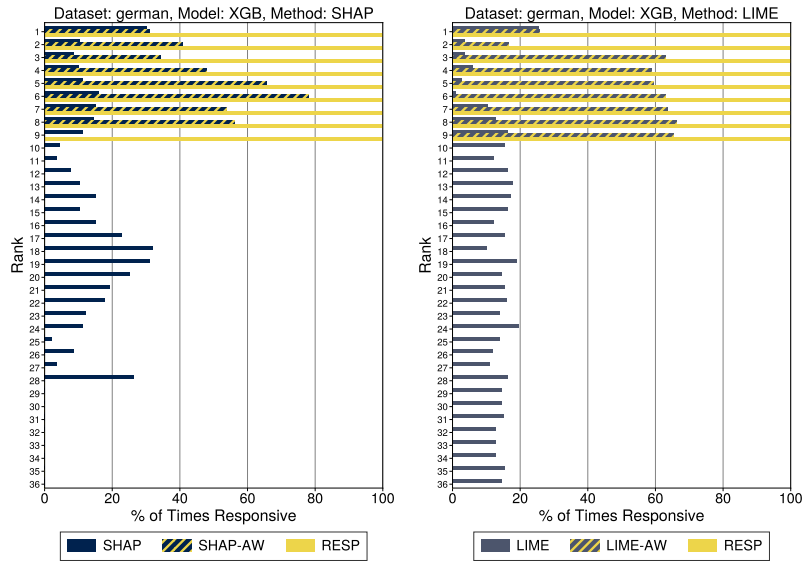
## C.2.2 `german`



**Figure 7:** Responsiveness of features for individuals who are denied credit by the LR model on the `german` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."
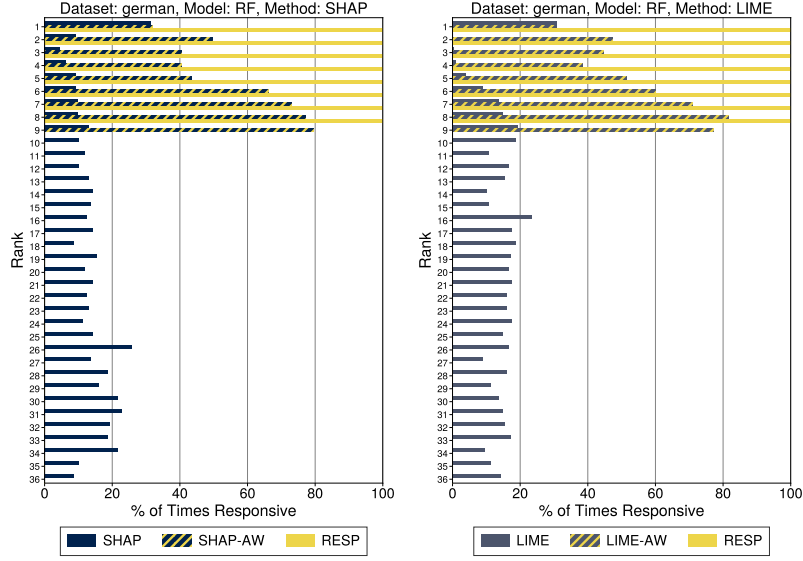


**Figure 8:** Responsiveness of features for individuals who are denied credit by the XGB model on the `german` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."

**Figure 9:** Responsiveness of features for individuals who are denied credit by the RF model on the `german` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."
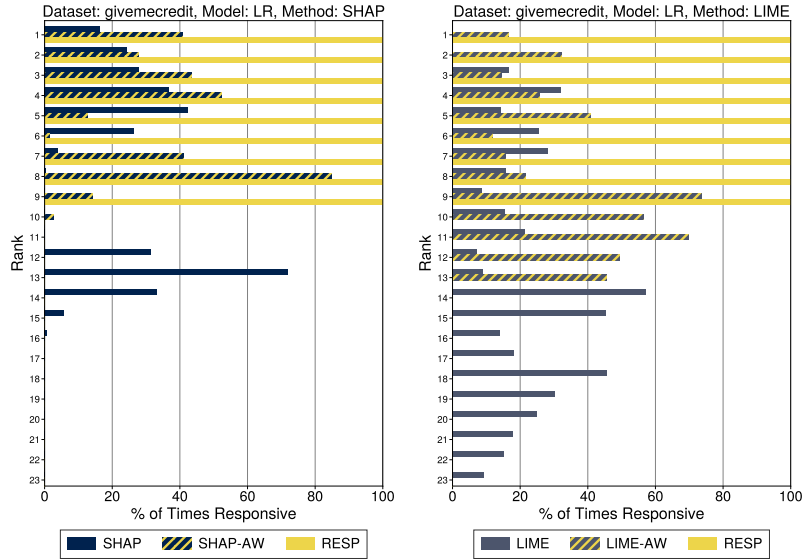
### C.2.3 `givemecredit`



**Figure 10:** Responsiveness of features for individuals who are denied credit by the LR model on the `givemecredit` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."
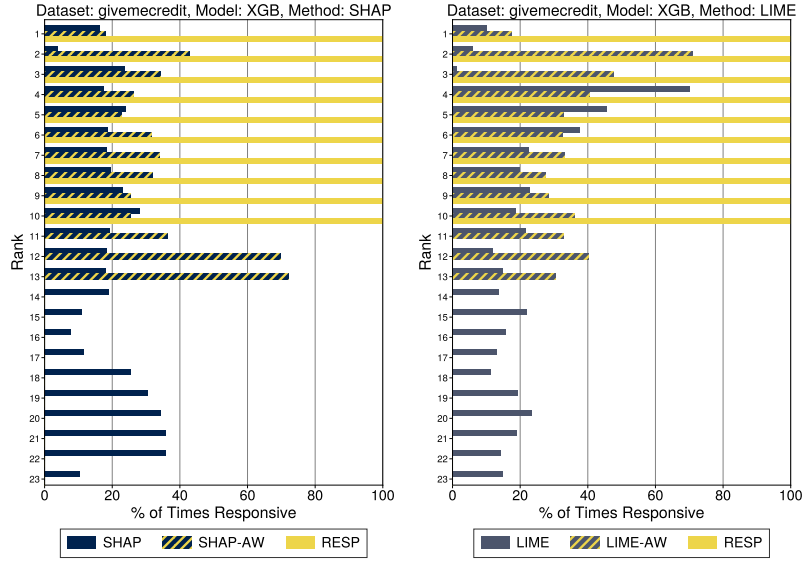
**Figure 11:** Responsiveness of features for individuals who are denied credit by the XGB model on the `givemecredit` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."
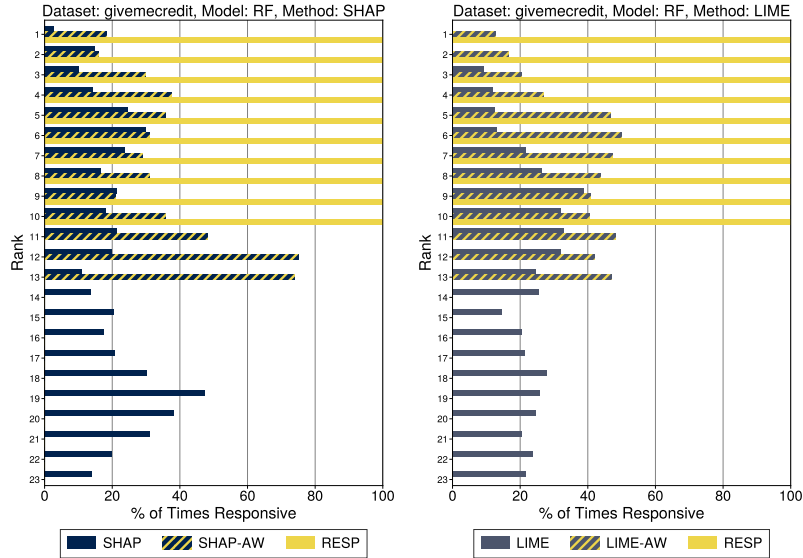


**Figure 12:** Responsiveness of features for individuals who are denied credit by the RF model on the `givemecredit` dataset according to absolute feature attribution rank using the original feature attribution method, its action-aware variant and RESP. For each method, we report the proportion of individuals with at least one responsive intervention on a feature with the $k$-th largest score ($k$-th ranked feature). Features must have non-zero score to be included in a "rank."

# D    SUPPLEMENTARY CASE STUDY DETAILS

## D.1    ACTIONABILITY CONSTRAINTS

The joint actionability constraints listed in Table 13 include:

| Name | Type | LB | UB | Actionability | Sign | Joint Constraints | Partition ID |
|---|---|---|---|---|---|---|---|
| Age | $\mathbb{Z}$ | 21 | 103 | No | | 0, 8, 10 | 0 |
| HistoryOfLatePaymentInPast2Years | $\{0,1\}$ | 0 | 1 | Yes | + | 0, 8, 10 | 0 |
| HistoryOfDelinquencyInPast2Years | $\{0,1\}$ | 0 | 1 | Yes | + | 0, 8, 10 | 0 |
| NumberRealEstateLoansOrLines | $\mathbb{Z}$ | 0 | 100 | Yes | + | 5, 6 | 5 |
| NumberOfOpenCreditLinesAndLoans | $\mathbb{Z}$ | 0 | 100 | Yes | + | 5, 6 | 5 |
| NumberOfDependents | $\mathbb{Z}$ | 0 | 20 | No | | – | 1 |
| DebtRatio | $\mathbb{R}$ | 0.0 | 61106.5 | Yes | | – | 2 |
| MonthlyIncome | $\mathbb{Z}$ | 0 | 3008750 | Yes | | – | 3 |
| CreditLineUtilization | $\mathbb{R}$ | 0.0 | 50708.0 | Yes | | – | 4 |
| HistoryOfLatePayment | $\{0,1\}$ | 0 | 1 | No | | – | 6 |
| HistoryOfDelinquency | $\{0,1\}$ | 0 | 1 | No | | – | 7 |

**Table 13:** Separable Actionability Constraints for the processed continuous `givemecredit` dataset. **Type** indicates the feature type ($\mathbb{Z}$ for integer, $\{0,1\}$ for binary). **LB**, **UB** are the lower and upper bounds for the feature. **Actionability** indicates whether the feature is globally actionable. **Sign** indicates if the feature can only increase (+) or decrease (-). **Joint Constraints** are a list non-separable constraint indices it is tied to (if any). **Partition ID** indicates which partition the feature belongs to.

1. DirectionalLinkage: Actions on `NumberRealEstateLoansOrLines` will induce to actions on `['NumberOfOpenCreditLinesAndLoans']`.Each unit change in `NumberRealEstateLoansOrLines` leads to:1.00-unit change in `NumberOfOpenCreditLinesAndLoans`

2. DirectionalLinkage: Actions on `HistoryOfLatePaymentInPast2Years` will induce to actions on `['Age']`.Each unit change in `HistoryOfLatePaymentInPast2Years` leads to:2.00-unit change in `Age`

3. DirectionalLinkage: Actions on `HistoryOfDelinquencyInPast2Years` will induce to actions on `['Age']`.Each unit change in `HistoryOfDelinquencyInPast2Years` leads to:2.00-unit change in `Age`

## D.2  MODEL PERFORMANCE

| | XGB | |
|---|---|---|
| Dataset | Train | Test |
| `givemecredit` $n = 120,268$ $d = 11$ Kaggle [35] | 0.937 | 0.830 |

**Table 14:** Model Performance of XGB model on the `givemecredit` dataset for Section 6