

A Data Description and Preprocessing

The datasets were divided into training and test sets using an 80-20 stratification. The numerical features were transformed by standardization to have a mean of zero and a variance of one. The categorical features, which have k different levels, were transformed into $k - 1$ binary variables using one-hot encoding. The binary characteristics were transformed into a single dummy variable using one-hot encoding. The sizes of the datasets before and after encoding are shown in Table 4.

	OBSERVATIONS	PRE-ENCODED FEATURES	POST-ENCODED FEATURES
COMPAS	6,907	7	7
Adult	32,561	14	107
MIMIC-III	48,786	14	14
Diabetes	101,766	33	101
German Credit	1,000	20	59
FICO	10,459	23	23
Headlines	41,752	12	17

Table 4: Training Dataset Sizes

Below we provide more details for each dataset.

COMPAS

The COMPAS dataset contains information on criminal recidivism in Broward County, Florida [Jeff Larson and Angwin, 2016]. The goal of this dataset is to predict the likelihood of recidivism within a two-year period, taking into account the following variables: gender, age, prior convictions, number of juvenile felonies/misdemeanors, and whether the current charge is a felony.

Adult

The Adult data is derived from U.S. Census statistics, including information on demographics, education, employment, marital status, and financial gain/loss [Dua and Graff, 2017]. The target variable of this dataset is whether an individual’s salary exceeds \$50,000.

MIMIC-III

MIMIC-III is a comprehensive database that stores a variety of medical data related to the experience of patients in the Intensive Care Unit (ICU) at Beth Israel Deaconess Medical Center [Johnson et al., 2016a,b]. The outcome of interest is determined by the binary indicator known as the “hospital expires flag,” which indicates whether or not a patient died during their hospitalization. We chose the following set of variables as features: age, preiculus (pre-ICU length of stay), gcs (Glasgow Coma Scale), heartrate_min, heartrate_max, meanbp_min (min blood pressure), meanbp_max (max blood pressure), resprate_min, resprate_max, tempc_min, tempc_max, urineoutput, mechvent (whether the patient is on mechanical ventilation), and electivesurgery (whether the patient had elective surgery).

Diabetes

The Diabetes dataset is derived from 10 years (1999-2008) of clinical care at 130 hospitals and integrated delivery networks in the United States [Dua and Graff, 2017]. It consists of more than 50 characteristics that describe patient and hospital outcomes. The dataset includes variables such as race, gender, age, admission type, time spent in hospital, specialty of admitting physician, number of lab tests performed, number of medications, and so on. We consider whether the patient will return to the hospital within 2 years as a binary indicator.

591 **German Credit**

592 The German credit data [Dua and Graff, 2017] uses financial and demographic indicators such
593 as checking account status, credit history, employment/marital status, etc., to predict whether an
594 individual will default on a loan.

595 **FICO**

596 The FICO Home Equity Line of Credit (HELOC) dataset [FICO, 2018] is used for the Explainable
597 Machine Learning Challenge. It includes a number of financial indicators, such as the number of
598 inquiries on a user’s account, the maximum delinquency, and the number of satisfactory transactions,
599 among others. These indicators relate to different individuals who have applied for credit. The target
600 variable is whether a consumer has been 90 or more days delinquent at any time within a 2-year
601 period since opening their account.

602 **Headlines**

603 The News Headline dataset [Chen et al., 2023b] is a survey data aimed at discovering what
604 kind of news content is shared and what factors are significantly associated with news shar-
605 ing. The survey includes several factors, including, age, income, gender, ethnicity, social
606 protection, economic protection, truth (“What is the likelihood that the above headline is
607 true?”), familiarity (“Are you familiar with the above headline (have you seen or heard about it
608 before?)? ”), Importance (“Assuming the headline is completely accurate, how important would
609 you consider this news to be?”), Political Concordance (“Assuming the above headline is com-
610 pletely accurate, how favorable would you consider it to be for Democrats versus Republicans?”).
611 The goal of this data set is to predict Sharing (“If you were to see the above article on social media,
612 how likely would you be to share it?”).

B Sensitivity of the reference points

In this section, we will mainly show how sensitive SEV^- is when we change the reference. Figure 8 shows an example of this, where moving the reference further away from the query (from r to the r') changes the SEV^- from 2 to 1. In this figure, the dark blue axes represent the feature values of different reference values, while the black dashed line represents the decision boundary of a linear classifier. Areas with different colors represent data points with different SEV^- . When the reference moves further from the decision boundary (from r to r'), the corresponding areas for SEV^- will move away from the decision boundary. For example, the star located in the yellow area has an SEV^- of 1 instead of 2 when the reference moves from r to r' . If the reference point is r , then the query needs to align the feature values along both x and y-axis to reach the SEV Explanation with reference r (recall an example of SEV^- explanation in Figure 2) in Section 3.2, which is the same point as r . However, if the reference point is r' , then the query only needs to align the feature value along the x-axis to reach the SEV Explanation with $SEV = 1$, which is the light blue dot.

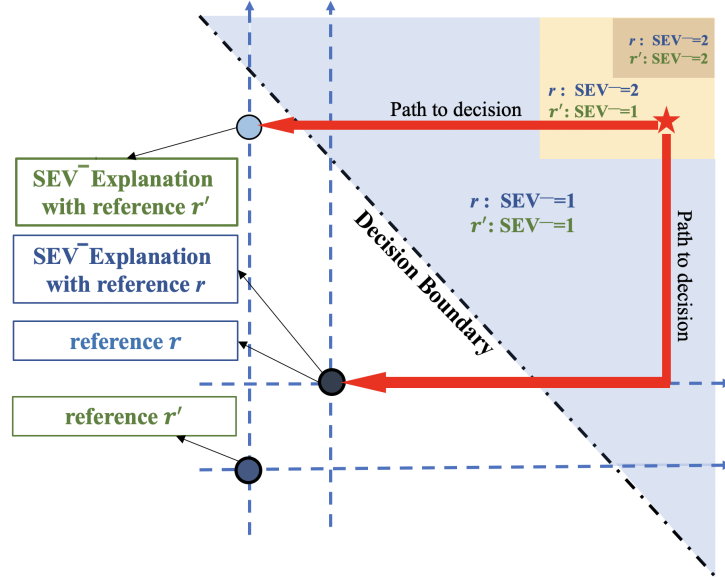


Figure 8: SEV^- distribution

Experiments have also shown that moving data points closer to the decision boundary might increase SEV^- . The result on the Explainable ML Challenge loan decision data [FICO, 2018] shown in Table 5 demonstrates that altering the reference point may increase the average SEV^- (from 3 to 5), but also introduces “unexplainable” samples (meaning $SEV^- \geq 10$). Hence, SEV^- is sensitive to the reference.

Table 5: SEV^- change by moving reference point \tilde{r} moving closer to the decision boundary to \tilde{r}'

MODEL	REFERENCE POINT	MEAN SEV^-	% OF SAMPLES		
			$SEV^- \geq 3$	$SEV^- \geq 6$	$SEV^- \geq 10$
L2LR	\tilde{r}	2.76	2.82	0	0
	\tilde{r}'	4.95	89.23	32.3	0
L1LR	\tilde{r}	2.46	1.00	0	0
	\tilde{r}'	4.57	56.87	21.27	0

C Detailed Description for Score-based Soft K-Means

As we have discussed in Section 4.1, SEV^- needs to have negatively predicted reference points. Therefore, when clustering the negative population, it is necessary to avoid positively predicted cluster centers. However, for most of the existing clustering methods, it is hard to “penalize” the positive predicted clusters, or their assigned samples. Therefore, we have modified the soft K-Means [Bezdek et al., 1984] algorithm so as to encourage negative clustering results.

The original Soft K-Means (SKM) algorithm generalizes K-means clustering by assigning membership scores for multiple clusters to each point. Given a data set $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ and C clusters, the goal is to minimize the objective function $J(U, V)$, where $U = [u_{ij}]$ is the membership matrix and $V = \{\mathbf{v}_1, \dots, \mathbf{v}_C\}$ are the weighted cluster centroids. The objective is to minimize:

$$J(U, V) = \sum_{i=1}^n \sum_{j=1}^C u_{ij}^m \|\mathbf{x}_i - \mathbf{v}_j\|_2^2 \quad (5)$$

where u_{ij} is the (soft) membership score of \mathbf{x}_i in cluster j :

$$u_{i,j} = \frac{1}{\sum_{k=1}^C \left(\frac{\|\mathbf{x}_i - \mathbf{v}_j\|_2}{\|\mathbf{x}_i - \mathbf{v}_k\|_2} \right)^{\frac{2}{m-1}}} \quad (6)$$

and $m > 1$ is a parameter that controls the strength towards each neighboring point. When $m \approx 1$, the SKM is similar to the performance of hard K-means clustering methods. When $m > 1$ for point i , it is considered to be associated with multiple clusters instead of one distinct cluster. The higher the value of m , the more a point is considered to be part of multiple clusters, thereby reducing the distinctness of each cluster and creating a more integrated and interconnected clustering arrangement. To avoid the cluster group being predicted positively, we have given higher m for those positive samples. Therefore, if the samples are predicted as positive, it reduces the possibility that those positively predicted samples to group as a cluster, which we can replace m as m'_i for each instance \mathbf{x}_i as

$$m'_i = 2m \cdot \min\{f(\mathbf{x}_i) - 0.5, 0\} + 1. \quad (7)$$

The value of $\min\{f(\mathbf{x}_i) - 0.5, 0\}$ increases as \mathbf{x}_i is classified as positive and further away from the decision boundary. As m' increases, the negatively predicted samples are more associated with one distinct cluster, while the positively predicted samples are associated with multiple clusters with smaller weight. This makes the cluster centers less likely to be influenced by positively predicted points. Thus, we can rewrite the objective of the soft K-Means algorithm can be modified as

$$J'(U, V) = \sum_{i=1}^n \sum_{j=1}^C u_{ij}^{m'_i} \|\mathbf{x}_i - \mathbf{v}_j\|_2^2. \quad (8)$$

We call this new objective function for encouraging negative clustering centers Score-based Soft K-Means (SSKM). In our experiments, the clustering is applied to the dataset after PaCMAP [Wang et al., 2021], and the feature mean of all samples in a cluster is considered as the cluster center of this cluster, which is eventually used as a reference point. The queries are assigned to reference points that are closest (based on ℓ_2 distance) to them in the PaCMAP embedding space for SEV° calculation. The reason why we would like to first embed the dataset is that the dimension of the datasets might be too high for direct clustering, and PaCMAP provides an embedding that preserves both local and global structure. Figure 9 shows the probability of the negative predicted instances, as well as the clustering results using different kinds of clustering methods. The red points and stars represent the positively predicted instances and cluster centers, while the blue ones are the negatively predicted instances and cluster centers. It is evident from the Figure that that SKM is more likely to introduce positively predicted cluster centers, compared to SSKM.

When we calculate SEV° in the experiments, all clustering parameters are tuned and fixed. For the rest of the datasets, the embedding using PaCMAP, and their clustering results for the negative population with their cluster centers, are shown in Figure 10. The regions with different colors represent different clusters, the blue stars in the graphs are cluster centers, and the gray points within the graphs are positive queries. All those cluster centers can be constrained to be predicted as negative by tuning the hyperparameter for Score-based Soft K-Means. Note that if one of the cluster centers cannot be constrained to be predicted as negative even with high m , then it is reasonable to remove this cluster center when calculating SEV° .

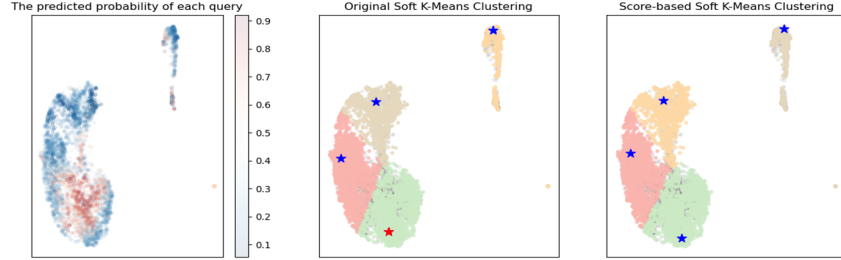


Figure 9: The clustering results for FICO dataset. (Left) The probability distribution for the negatively labeled queries; (Middle) The clustering result for Original Soft K-Means Clustering; (Right) The clustering result for Score-based K-Means Clustering. The red stars represent the positively predicted cluster centers, and the blue stars the negatively predicted cluster centers.

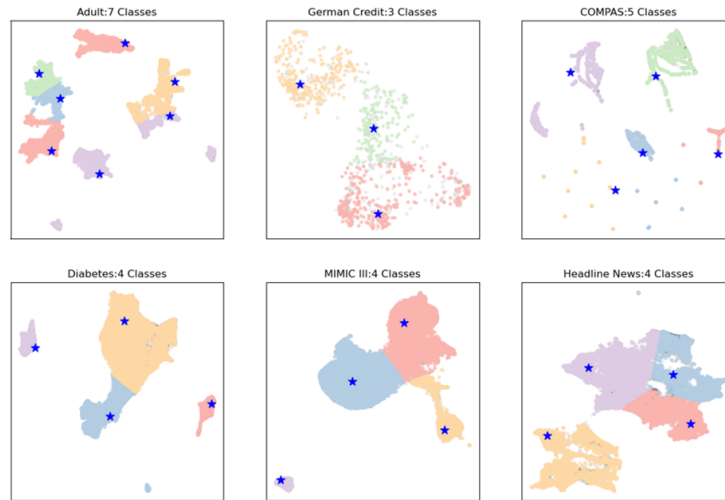


Figure 10: Clustering Results for different datasets.

676 D Detailed Algorithm for Flexible-based SEV

677 This section presents how the flexible-based SEV (SEV^F) has done to determine the flexible refer-
 678 ences. The key idea of finding the reference is to do a grid search through each of the features in the
 679 training dataset based on the original reference, and find the feature values that has the minimum
 680 model outcome.

Algorithm 1 Reference Search for Flexible SEV

```

1: Input: The negative samples  $X^-$ , flexibility  $\epsilon$ , reference  $\tilde{r}$ , grid size  $G$ 
2: Output: Flexible reference  $\tilde{r}'$ 
3: Initialization:  $\tilde{r}' \leftarrow \tilde{r}$ 
4: for each feature  $j \in \mathcal{J}$ , where  $\tilde{r}_j$  is the reference value of feature  $j$  in  $X^-$  do
5:    $q_j \leftarrow \text{quantile}(X_j^-, \tilde{r}_j)$  {Quantile location of  $\tilde{r}_j$ }
6:    $B_j^+ \leftarrow \text{percentile}(X_j^-, q_j + \epsilon)$  {The upper range}
7:    $B_j^- \leftarrow \text{percentile}(X_j^-, q_j - \epsilon)$  {The lower range}
8:    $B_j^{(g)} \sim \text{Uniform}[B_j^-, B_j^+], g = 1, \dots, G$ 
9:    $P_j^{(g)} \leftarrow f([\tilde{r}_1, \dots, B_j^{(g)}, \dots, \tilde{r}_J]), g = 1 \dots G$  {Slight change to feature  $j$  for prediction}
10:   $g' \leftarrow \arg \min_g P_j^{(g)}$  {Find minimum model outcome}
11:   $\tilde{r}'_j \leftarrow B_j^{(g')}$  {Update for flexible references}
12: end for

```

E Detailed Algorithms for Tree-based SEV

This section presents how the tree-based SEV is calculated through two main procedure: Algorithm 2 (Preprocessing) for collecting all negative pathways and assigning them to each internal nodes and Algorithm 3 (Efficient SEV^T Calculation) for checking all negative pathways conditions for each query and calculating the number of feature changes.

Algorithm 2 Preprocessing - Information collection process for SEV^T

```

1: Input: Decision tree  $DT$ 
2: Output:  $DT^-$ , a dictionary of paths to negative predictions for each internal node encoding
3:  $nodes \leftarrow [DT.root]$ 
4:  $negative\_path \leftarrow []$ 
5: {Negative path collection procedure}
6: while  $nodes$  not empty do
7:    $[node, path] \leftarrow nodes.pop()$ 
8:   if  $node$  is a negative leaf then
9:      $negative\_path.append(path)$ 
10:  else if  $node$  is an internal node or a root node then
11:    {Add the child nodes and the path to the node list}
12:     $nodes.append([node.left, path + "L"])$ 
13:     $nodes.append([node.right, path + "R"])$ 
14:  else
15:    Continue {if the leaf is positive, ignore it}
16:  end if
17: end while
18: {Assign Negative Pathways to root or internal nodes}
19:  $DT^- \leftarrow dict()$ 
20: for each  $path$  in  $negative\_path$  do
21:   for  $i = 1, \dots, path.length$  do
22:    {Add the negative decision path for internal nodes}
23:     $curr\_node \leftarrow negative\_path[:i]$ 
24:    { $curr\_node$  is the encoded internal node, and  $negative\_path[i:]$  is a negative decision path below this node}
25:     $DT^-[curr\_node].append(negative\_path[i:])$ 
26:   end for
27: end for

```

Algorithm 3 Efficient SEV^T Calculation – Negative Pathways Check

```
1: Input:  $DT$ : decision tree,  $DT^-$ : decision trees with paths to negative predictions, query value  $x_i$ ,  $DP_i$ : list of internal nodes representing decision process for  $x_i$ ,  $path_i$ : the encoded  $DP_i$ 
2: Output:  $SEV^T$ 
3: INITIALIZATION:  $SEV^T \leftarrow 0$ 
4:  $decision\_path \leftarrow encoded(DT, x_i)$ 
5: { $encoded(DT, x_i)$  is a function to get the string representation of the query  $x_i$  or a node  $node$  for  $DT$ , e.g. "LR", "LL" mentioned in section 4.2}
6: for each internal node  $node$  in  $DP_i$  do
7:   if  $node$  has a sibling leaf node and is predicted as negative then
8:      $SEV^T \leftarrow 1$  {Based on Theorem 4.1}
9:     Break { $SEV^T=1$  is the smallest  $SEV^T$ , no further calculation needed}
10:  end if
11:   $encoded\_node \leftarrow encoded(DT, node)$  {Get the string representation of  $node$ }
12:   $negative\_paths \leftarrow DT^-[encoded\_node]$  {Get the negative pathways  $encoded\_node$  have}
13:  for each  $path$  in  $negative\_paths$  do
14:    {If the negative goes the same direction as the decision path, we don't need to calculate this path again}
15:    { $path[0]$  is the first character in  $path$ }
16:    if  $decision\_path[encoded\_node.length]=path[0]$  then
17:      Continue
18:    end if
19:     $temp\_sev \leftarrow 0$ 
20:    {Go over the condition in the  $path$ }
21:    {Check if query  $x_i$  satisfies, if it doesn't satisfy the condition, then  $temp\_sev$  should add 1}
22:    for  $condition$  in each  $path$  do
23:      if  $x_i$  doesn't satisfy  $condition$  then
24:         $temp\_sev \leftarrow temp\_sev + 1$ 
25:      end if
26:    end for
27:     $SEV^T \leftarrow \min\{temp\_sev, SEV^T\}$  {Update  $SEV^T$  to be the smaller one}
28:    if  $SEV^T = 1$  then
29:      Break { $SEV^T=1$  is the smallest  $SEV^T$ , no further calculation needed}
30:    end if
31:  end for
32: end for
```

F Model Training and parameters selection

Baseline models were fit using `sklearn` [Pedregosa et al., 2011] implementations in Python. The logistic regression models L1 LR and L2 LR were fit using regularization parameter $C = 0.01$. The 2-layer MLP used ReLU activation and consisted of two fully-connected layers with 128 nodes each. It was trained with early stopping. The gradient-boosted classifier used 200 trees with a max depth of 3. For tree-based methods comparisons, the decision tree classifiers were fit using `sklearn` [Pedregosa et al., 2011] and `TreeFARMS` packages [Wang et al., 2022b]. Since GOSDT methods require binary input, we used the built-in threshold guessing function in GOSDT to binarize the features with set of parameters `n_est=50`, and `max_depth=1`. All the models are trained using a RTX2080Ti GPU, and with 4 core in Intel(R) Xeon(R) Gold 6226 CPU @ 2.70GHz.

In order to test the performance of All-Opt⁻, all models mentioned above were trained by adding the SEV losses from Section 5 to the standard loss term (BCELoss). For GBDT, the training goal is to reweigh the trees from the baseline GBDT model. The resulting loss was minimized via gradient descent in PyTorch [Paszke et al., 2019], with a batch size of 128, a learning rate of 0.1, and the Adam optimizer. To maintain high accuracy, the first 80 training epochs are warm-up epochs optimizing just Binary Cross Entropy Loss for classification (BCELoss). The next 20 epochs add the All-Opt terms and the baseline positive penalty term to encourage low SEV values. Moreover, during the optimization process, it is important to ensure that the reference has a negative prediction. If the reference is predicted as positive, then the SEV⁻ may not exist, and a sparse explanation is no longer meaningful. Thus, we add a term to penalize the reference if it receives a positive prediction:

$$\ell_{\text{Pos_ref}}(f) := \sum_{i=1}^n \max(f(\tilde{\mathbf{r}}_i), 0.5 - \theta)$$

where $\theta > 0$ is a margin parameter, usually $\theta = 0.05$. This term is $(0.5 - \theta)$ as long as the reference is predicted negative. As soon as it exceeds that amount, it is penalized (increasing linearly in $f(\tilde{\mathbf{r}})$).

To put these into an algorithm, we optimize a linear combination of different loss terms,

$$\min_{f \in \mathcal{F}} \ell_{\text{BCE}}(f) + C_1 \ell_{\text{SEV_All_Opt}^-}(f) + C_2 \ell_{\text{Pos_ref}}(f) \quad (9)$$

Therefore, we are tuning both C_1 and C_2 to find a model with sparser explanations without performance loss through grid search. For cluster-based SEV, the cluster centers are recalculated based on the new model every 5 epochs.

G The sparsity and meaningful performance of different counterfactual explanation methods

In this section, we provide detailed information on other kinds of counterfactual explanations generated by the CARLA package [Pawelczyk et al., 2021] on different datasets for logistic regression models. Table 6 shows the number of features changed and the ℓ_∞ for different counterfactual explanations. These counterfactual explanations tend to provide less sparse explanations than other SEV^- variants shown in Section 6.3. For the ℓ_∞ calculations, we consider only the numerical features, since the categorical features' ℓ_∞ norm does not provide meaningful explanations. Moreover, we have calculated the average log-likelihood of the explanations using the Gaussian Mixture Model in scikit-learn Pedregosa et al. [2011]. The parameter `n_components` for each dataset is selected based on the clustering result mentioned in Appendix C. Here, we are using the same Gaussian Mixture Model for evaluating whether the explanation is within a high-density region.

Table 6: Explanation performance in different counterfactual explanations

DATASET	COUNTERFACTUAL EXPLANATIONS	MEAN ℓ_∞	# FEATURES CHANGE	MEDIAN LOG-LIKELIHOOD
Adult	Growing Sphere	1.07 ± 0.01	14 ± 0.00	345.03 ± 34.19
	DiCE	0.78 ± 0.02	2.19 ± 0.12	-24752.12 ± 452.47
	REVISE	6.1 ± 0.02	12.14 ± 0.75	345.03 ± 32.84
	Watcher	0.01 ± 0.01	6.00 ± 0.00	345.12 ± 34.19
	SEV^1	22.62 ± 0.01	1.18 ± 0.02	-24752.12 ± 452.47
	SEV°	2.86 ± 0.01	1.34 ± 0.02	156.88 ± 59.67
COMPAS	Growing Sphere	0.02 ± 0.01	7.00 ± 0.00	10.47 ± 0.00
	DiCE	1.38 ± 0.02	3.20 ± 0.45	-6.68 ± 0.02
	REVISE	1.12 ± 0.03	5.54 ± 0.63	-1.84 ± 0.21
	Watcher	0.01 ± 0.01	5.00 ± 0.00	10.48 ± 0.03
	SEV^1	2.31 ± 0.01	1.22 ± 0.02	14.65 ± 0.32
	SEV°	2.06 ± 0.01	1.19 ± 0.02	14.41 ± 0.05
Diabetes	Growing Sphere	0.01 ± 0.01	33.00 ± 0.00	320.41 ± 21.47
	DiCE	0.71 ± 0.12	2.76 ± 0.15	-74296.98 ± 861.27
	REVISE	0.80 ± 0.02	15.84 ± 0.02	320.41 ± 16.73
	Watcher	0.01 ± 0.01	12 ± 0.00	320.41 ± 21.34
	SEV^1	2.7 ± 0.10	1.63 ± 0.01	309.56 ± 15.32
	SEV°	2.31 ± 0.12	1.28 ± 0.02	320.71 ± 14.79
FICO	Growing Sphere	0.01 ± 0.01	23 ± 0.00	-10.93 ± 0.42
	DiCE	1.15 ± 0.13	3.27 ± 0.17	-20.11 ± 0.3
	REVISE	0.12 ± 0.01	23 ± 0.00	-10.94 ± 0.42
	Watcher	0.01 ± 0.01	23 ± 0.00	-10.94 ± 0.41
	SEV^1	1.81 ± 0.01	2.76 ± 0.02	-20.11 ± 0.32
	SEV°	1.82 ± 0.01	2.21 ± 0.02	-19.32 ± 0.21
German Credit	Growing Sphere	0.01 ± 0.02	20 ± 0.00	52.20 ± 0.02
	DiCE	6.08 ± 0.01	2.76 ± 0.23	-53908.78 ± 367.84
	REVISE	0.16 ± 0.01	7.65 ± 0.12	-73492.06 ± 492.45
	Watcher	0.01 ± 0.00	6.00 ± 0.00	52.23 ± 0.04
	SEV^1	3.08 ± 0.01	1.51 ± 0.02	-124914.32 ± 792.52
	SEV°	3.2 ± 0.01	1.17 ± 0.02	50.21 ± 0.32
Headline	Growing Sphere	0.01 ± 0.00	18 ± 0.00	-4.56 ± 0.02
	DiCE	1.13 ± 0.02	2.79 ± 0.14	-12.84 ± 0.42
	REVISE	1.81 ± 0.13	15.93 ± 0.24	-6.98 ± 0.12
	Watcher	0.01 ± 0.01	12 ± 0.00	-4.56 ± 0.02
	SEV^1	2.50 ± 0.02	1.98 ± 0.01	1.52 ± 0.12
	SEV°	2.94 ± 0.02	1.62 ± 0.02	0.89 ± 0.26
MIMIC	Growing Sphere	0.01 ± 0.01	14 ± 0.00	-24.52 ± 0.02
	DiCE	1.34 ± 0.23	6.47 ± 0.24	-26.55 ± 0.02
	REVISE	0.01 ± 0.00	12 ± 0.00	-24.52 ± 0.01
	Watcher	0.01 ± 0.00	12 ± 0.00	-24.52 ± 0.01
	SEV^1	4.53 ± 0.49	1.18 ± 0.02	-20.11 ± 0.32
	SEV°	1.98 ± 0.13	1.19 ± 0.02	-19.32 ± 0.15

H Detailed SEV^- for all datasets

In this section, we show how SEV^1 , SEV° , $SEV^{\circ+F}$ can increase the similarity metrics or reduce the sparsity explanations. All the models are trained and evaluated 10 times using different splits, and evaluated for their mean SEV^- , mean ℓ_∞ , as well as their explanation time for each query.

Table 7 shows the model performance and SEV^1 on various datasets. SEV^1 is considered as a base case for other SEV^- variants to compare with. Table 7 shows that SEV^1 yields very high ℓ_∞ for each model, indicating a large distance between the query and reference, which implies low closeness according to Section 3.2.

Table 8 shows the model performance and SEV° on different datasets. Similarly, The Mean SEV° column reports the mean SEV° for the model and the decrease in mean SEV^- in percentage compared to SEV^1 (reported in the parenthesis). The Mean ℓ_∞ column reports the mean ℓ_∞ and the percentage reduction compared to SEV^1 . On most datasets, SEV° increases, and ℓ_∞ decreases, which means that the model is providing both sparser and more meaningful explanations. For some datasets like Adult and MIMIC, the SEV° increases, since the cluster-based reference points might be closer to the decision boundary of the model as each query is trying to find the closest (in ℓ_2 distance) negatively predicted reference point, which might provide less sparse explanations.

Table 9 shows the model performance and $SEV^{\circ+F}$ (SEV° with variable reference) on various datasets with different flexibility levels. The Mean SEV^F column reports the mean SEV^- for the model and the decrease in mean SEV^- in percentage compared to SEV^1 (reported in the parenthesis). The Mean ℓ_∞ column reports the mean ℓ_∞ and the percentage reduction compared to SEV^1 . It is evident that with SEV^F , SEV^- decreases, but the ℓ_∞ norm will increase due to the flexibility of the features mentioned in section 4.4. The “flexibility used” column shows the proportion of queries using the flexible reference instead of the original one for calculating SEV^F , and the higher the proportion, the larger decrease in SEV^- the model can achieve.

Table 7: The SEV^1 under different models

DATASET	MODEL	TRAIN ACCURACY	TEST ACCURACY	TRAIN AUC	TEST AUC	AVERAGE SEV^1	MEDIAN ℓ_∞	EXPLANATION TIME(10^{-2} s)	AVERAGE LOG-LIKELIHOOD
Adult	GBDT	0.88 \pm 0.0	0.87 \pm 0.0	0.93 \pm 0.0	0.93 \pm 0.0	1.23 \pm 0.02	18.28 \pm 1.8	0.69 \pm 0.08	-57437.86 \pm 2718.7
	L1LR	0.85 \pm 0.0	0.85 \pm 0.0	0.9 \pm 0.0	0.9 \pm 0.0	1.14 \pm 0.01	24.2 \pm 2.41	0.26 \pm 0.01	-44735.07 \pm 1393.91
	L2LR	0.85 \pm 0.0	0.85 \pm 0.0	0.9 \pm 0.0	0.9 \pm 0.0	1.18 \pm 0.0	22.62 \pm 2.27	0.16 \pm 0.01	-49293.12 \pm 1157.19
	MLP	0.87 \pm 0.0	0.86 \pm 0.0	0.93 \pm 0.0	0.92 \pm 0.0	1.27 \pm 0.06	21.73 \pm 3.57	0.62 \pm 0.17	-67000.48 \pm 5030.26
COMPAS	GBDT	0.7 \pm 0.0	0.67 \pm 0.01	0.77 \pm 0.0	0.72 \pm 0.01	1.15 \pm 0.04	1.94 \pm 0.08	0.18 \pm 0.02	8.15 \pm 0.97
	L1LR	0.68 \pm 0.0	0.67 \pm 0.01	0.73 \pm 0.0	0.72 \pm 0.01	1.25 \pm 0.02	2.31 \pm 0.07	0.12 \pm 0.0	5.09 \pm 0.92
	L2LR	0.68 \pm 0.0	0.67 \pm 0.02	0.73 \pm 0.0	0.72 \pm 0.01	1.26 \pm 0.03	2.41 \pm 0.09	0.08 \pm 0.01	5.19 \pm 1.0
	MLP	0.69 \pm 0.01	0.67 \pm 0.01	0.74 \pm 0.01	0.72 \pm 0.01	1.35 \pm 0.12	2.3 \pm 0.32	0.27 \pm 0.09	6.49 \pm 1.1
Diabetes	GBDT	0.65 \pm 0.0	0.64 \pm 0.0	0.66 \pm 0.0	0.66 \pm 0.0	1.39 \pm 0.01	2.82 \pm 0.01	364.74 \pm 92.38	-59814.81 \pm 2356.74
	L1LR	0.62 \pm 0.0	0.62 \pm 0.0	0.66 \pm 0.0	0.66 \pm 0.0	1.62 \pm 0.01	2.6 \pm 0.01	106.63 \pm 79.76	-20834.12 \pm 1378.32
	L2LR	0.62 \pm 0.0	0.62 \pm 0.0	0.66 \pm 0.0	0.66 \pm 0.0	1.63 \pm 0.01	2.7 \pm 0.01	117.63 \pm 79.76	-19117.45 \pm 1091.56
	MLP	0.65 \pm 0.01	0.64 \pm 0.0	0.71 \pm 0.01	0.69 \pm 0.0	1.69 \pm 0.13	2.67 \pm 0.09	136.33 \pm 140.47	-70595.3 \pm 3666.52
FICO	GBDT	0.71 \pm 0.0	0.7 \pm 0.0	0.78 \pm 0.0	0.77 \pm 0.01	3.58 \pm 0.12	1.81 \pm 0.01	692.83 \pm 30.77	-74.13 \pm 8.92
	L1LR	0.71 \pm 0.0	0.7 \pm 0.0	0.78 \pm 0.0	0.77 \pm 0.01	2.47 \pm 0.11	1.81 \pm 0.07	100.83 \pm 30.77	-81.31 \pm 7.41
	L2LR	0.72 \pm 0.0	0.71 \pm 0.01	0.78 \pm 0.0	0.78 \pm 0.01	2.76 \pm 0.12	1.93 \pm 0.04	481.75 \pm 146.53	-52.09 \pm 2.1
	MLP	0.72 \pm 0.01	0.71 \pm 0.01	0.8 \pm 0.02	0.78 \pm 0.01	2.7 \pm 0.29	1.88 \pm 0.15	553.15 \pm 463.34	-67.71 \pm 13.05
German Credit	GBDT	0.96 \pm 0.01	0.75 \pm 0.02	0.99 \pm 0.0	0.77 \pm 0.02	1.39 \pm 0.12	1.87 \pm 0.46	2.69 \pm 1.8	-75811.5 \pm 6476.74
	L1LR	0.75 \pm 0.01	0.75 \pm 0.01	0.8 \pm 0.01	0.79 \pm 0.05	1.3 \pm 0.06	2.45 \pm 0.16	0.78 \pm 0.49	-64237.32 \pm 26906.43
	L2LR	0.78 \pm 0.01	0.76 \pm 0.03	0.83 \pm 0.01	0.79 \pm 0.04	1.51 \pm 0.15	3.08 \pm 0.42	1.34 \pm 0.96	-111945.26 \pm 9916.8
	MLP	0.81 \pm 0.04	0.76 \pm 0.03	0.87 \pm 0.04	0.78 \pm 0.04	1.6 \pm 0.19	2.69 \pm 0.45	7.68 \pm 5.59	-119557.08 \pm 15328.57
Headline	GBDT	0.82 \pm 0.0	0.81 \pm 0.0	0.9 \pm 0.0	0.89 \pm 0.0	1.82 \pm 0.03	2.35 \pm 0.02	16.25 \pm 2.45	-395.41 \pm 340.77
	L1LR	0.78 \pm 0.0	0.78 \pm 0.0	0.85 \pm 0.0	0.85 \pm 0.0	1.92 \pm 0.01	2.51 \pm 0.02	6.73 \pm 0.38	-558.81 \pm 287.68
	L2LR	0.78 \pm 0.0	0.78 \pm 0.0	0.86 \pm 0.0	0.85 \pm 0.0	1.98 \pm 0.01	2.5 \pm 0.02	9.21 \pm 0.49	-555.95 \pm 286.15
	MLP	0.83 \pm 0.01	0.81 \pm 0.0	0.91 \pm 0.01	0.89 \pm 0.0	2.03 \pm 0.03	2.31 \pm 0.07	26.25 \pm 2.45	-493.37 \pm 316.22
MIMIC	GBDT	0.91 \pm 0.0	0.9 \pm 0.0	0.87 \pm 0.0	0.85 \pm 0.0	1.18 \pm 0.02	1.28 \pm 0.15	1.03 \pm 0.22	-18.92 \pm 0.37
	L1LR	0.89 \pm 0.0	0.89 \pm 0.0	0.8 \pm 0.0	0.8 \pm 0.0	1.15 \pm 0.02	4.53 \pm 0.49	0.26 \pm 0.04	-19.76 \pm 0.52
	L2LR	0.89 \pm 0.0	0.89 \pm 0.0	0.8 \pm 0.0	0.8 \pm 0.0	1.16 \pm 0.02	4.34 \pm 0.52	0.29 \pm 0.03	-19.66 \pm 0.49
	MLP	0.9 \pm 0.0	0.9 \pm 0.0	0.87 \pm 0.01	0.85 \pm 0.0	1.18 \pm 0.03	2.08 \pm 0.35	0.79 \pm 0.19	-17.25 \pm 0.84

Table 8: The SEV[®] under different models

DATASET	MODEL	TRAIN ACCURACY	TEST ACCURACY	TRAIN AUC	TEST AUC	AVERAGE SEV	MEDIAN ℓ_∞	AVERAGE TIME (10^{-2})	AVERAGE LOG- LIKELIHOOD
Adult	GBDT	0.88 ± 0.0	0.87 ± 0.0	0.93 ± 0.0	0.93 ± 0.0	1.39(13.01%)	2.41(-86.82%)	2.22 ± 0.84	-22974.51(60.0%)
	L1LR	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.23(7.89%)	2.05(-91.53%)	0.56 ± 0.03	-39333.37(12.07%)
	L2LR	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.34(13.56%)	2.86(-87.36%)	0.38 ± 0.12	-21033.54(57.33%)
	MLP	0.87 ± 0.0	0.86 ± 0.0	0.93 ± 0.0	0.92 ± 0.0	1.62(27.56%)	5.16(-76.25%)	1.18 ± 0.53	-23421.5(60.97%)
COMPAS	GBDT	0.77 ± 0.0	0.67 ± 0.01	0.77 ± 0.0	0.72 ± 0.01	1.18(2.61%)	1.52(-21.65%)	0.32 ± 0.03	9.08(11.41%)
	L1LR	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.19(-4.8%)	1.75(-24.24%)	0.12 ± 0.01	5.53(8.64%)
	L2LR	0.68 ± 0.0	0.67 ± 0.02	0.73 ± 0.0	0.72 ± 0.01	1.22(-3.17%)	2.06(-14.52%)	0.09 ± 0.01	5.98(15.22%)
	MLP	0.69 ± 0.01	0.67 ± 0.01	0.74 ± 0.01	0.72 ± 0.01	1.3(-3.7%)	1.82(-20.87%)	0.15 ± 0.03	9.12(40.52%)
Diabetes	GBDT	0.65 ± 0.0	0.64 ± 0.0	0.7 ± 0.0	0.7 ± 0.0	1.36(-2.21%)	1.89(-49.21%)	17.39 ± 7.21	-5572.49(90.55%)
	L1LR	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.22(-24.6%)	2.31(-11.58%)	2.1 ± 0.4	-5460.38(92.27%)
	L2LR	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.28(-21.47%)	2.31(-14.44%)	3.8 ± 1.26	-14461.36(24.36%)
	MLP	0.65 ± 0.0	0.63 ± 0.0	0.7 ± 0.01	0.69 ± 0.0	1.47(-13.02%)	2.24(-16.1%)	23.28 ± 14.31	-11320.72(83.96%)
FICO	GBDT	0.77 ± 0.0	0.72 ± 0.01	0.85 ± 0.0	0.79 ± 0.01	2.06(-42.52%)	1.08(-40.3%)	23.34 ± 8.86	-59.52(19.7%)
	L1LR	0.71 ± 0.0	0.7 ± 0.0	0.78 ± 0.0	0.77 ± 0.0	1.79(-27.53%)	1.95(7.73%)	3.11 ± 1.02	-77.53(4.65%)
	L2LR	0.72 ± 0.0	0.71 ± 0.01	0.78 ± 0.0	0.77 ± 0.01	2.21(-19.93%)	1.82(-5.7%)	39.49 ± 16.49	-58.86(-13.0%)
	MLP	0.74 ± 0.01	0.71 ± 0.01	0.81 ± 0.01	0.78 ± 0.01	2.15(-20.37%)	1.75(-6.91%)	26.26 ± 9.01	-62.6(7.55%)
German Credit	GBDT	0.96 ± 0.01	0.75 ± 0.02	0.99 ± 0.0	0.77 ± 0.03	1.22(-12.23%)	1.73(-7.49%)	0.79 ± 0.53	-28478.65(62.43%)
	L1LR	0.75 ± 0.01	0.75 ± 0.02	0.8 ± 0.01	0.77 ± 0.04	1.03(-20.77%)	1.52(-37.96%)	0.05 ± 0.01	-23691.73(63.12%)
	L2LR	0.78 ± 0.01	0.76 ± 0.03	0.83 ± 0.01	0.79 ± 0.04	1.17(-22.52%)	3.2(3.9%)	0.1 ± 0.07	-40622.35(63.71%)
	MLP	0.81 ± 0.04	0.76 ± 0.03	0.87 ± 0.04	0.78 ± 0.04	1.24(-22.5%)	2.54(-5.58%)	0.24 ± 0.2	-40045.69(66.5%)
Headline	GBDT	0.82 ± 0.0	0.81 ± 0.0	0.9 ± 0.0	0.89 ± 0.0	1.76(-3.3%)	2.18(-7.23%)	6.96 ± 0.84	-383.24(-3.08%)
	L1LR	0.78 ± 0.0	0.78 ± 0.0	0.85 ± 0.0	0.85 ± 0.0	1.57(-18.23%)	2.94(17.13%)	0.88 ± 0.21	-559.35(0.1%)
	L2LR	0.78 ± 0.0	0.78 ± 0.0	0.86 ± 0.0	0.85 ± 0.0	1.62(-18.18%)	2.94(17.6%)	1.46 ± 0.1	-556.52(0.1%)
	MLP	0.83 ± 0.01	0.81 ± 0.0	0.91 ± 0.01	0.89 ± 0.0	1.67(-17.7%)	1.99(-16.08%)	3.05 ± 0.43	-495.08(0.0%)
MIMIC	GBDT	0.91 ± 0.0	0.9 ± 0.0	0.87 ± 0.0	0.85 ± 0.0	1.21(2.54%)	0.49(-61.72%)	0.61 ± 0.12	-18.15(4.07%)
	L1LR	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.17(1.74%)	1.8(-60.26%)	0.17 ± 0.03	-20.41(-3.29%)
	L2LR	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.19(2.59%)	1.98(-54.38%)	0.19 ± 0.03	-20.26(-3.05%)
	MLP	0.9 ± 0.0	0.9 ± 0.0	0.87 ± 0.01	0.85 ± 0.0	1.23(4.24%)	0.6(-71.15%)	0.33 ± 0.07	-16.77(2.78%)

Table 9: SEV⁺_F under different models

DATASET	MODEL	FLEX- IBILITY	TRAIN ACCURACY	TEST ACCURACY	TRAIN AUC	TEST AUC	AVERAGE SEV ⁻	MEDIAN ℓ_∞	AVERAGE LOG- LIKELIHOOD	EXPLANATION TIME(10^{-2} s)
Adult	GBDT	0.05	0.88 ± 0.0	0.87 ± 0.0	0.93 ± 0.0	0.93 ± 0.0	1.3(5.69%)	0.95(-94.8%)	-21763.14(62.11%)	3.98 ± 0.45
		0.10	0.88 ± 0.0	0.87 ± 0.0	0.93 ± 0.0	0.93 ± 0.0	1.29(4.88%)	0.95(-94.8%)	-20395.38(4.49%)	3.82 ± 0.32
		0.20	0.88 ± 0.0	0.87 ± 0.0	0.93 ± 0.0	0.93 ± 0.0	1.29(4.88%)	0.96(-94.75%)	-17611.65(69.34%)	3.63 ± 0.29
	L1LR	0.05	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.2(5.26%)	0.96(-96.03%)	-29801.44(33.38%)	1.0 ± 0.04
		0.10	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.19(4.39%)	0.96(-96.03%)	-29144.93(34.85%)	0.94 ± 0.04
		0.20	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.19(4.39%)	0.97(-95.99%)	-30245.09(32.39%)	0.91 ± 0.04
	L2LR	0.05	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.32(11.86%)	2.47(-89.08%)	-20693.31(58.02%)	1.59 ± 0.19
		0.10	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.32(11.86%)	2.41(-89.35%)	-20294.61(58.83%)	1.64 ± 0.18
		0.20	0.85 ± 0.0	0.85 ± 0.0	0.9 ± 0.0	0.9 ± 0.0	1.32(11.86%)	2.49(-88.99%)	-21987.43(55.39%)	1.59 ± 0.16
	MLP	0.05	0.87 ± 0.0	0.86 ± 0.0	0.93 ± 0.0	0.92 ± 0.0	1.54(21.26%)	2.95(-86.42%)	-27141.97(59.49%)	3.78 ± 1.4
		0.10	0.87 ± 0.0	0.86 ± 0.0	0.93 ± 0.0	0.92 ± 0.0	1.52(19.69%)	2.75(-87.34%)	-23444.97(65.01%)	3.76 ± 1.36
		0.20	0.87 ± 0.0	0.86 ± 0.0	0.93 ± 0.0	0.92 ± 0.0	1.44(13.39%)	2.37(-89.09%)	-22225.46(66.83%)	2.88 ± 1.11
COMPAS	GBDT	0.05	0.7 ± 0.0	0.67 ± 0.01	0.77 ± 0.0	0.72 ± 0.01	1.2(4.35%)	1.44(-25.77%)	8.85(8.59%)	0.77 ± 0.06
		0.10	0.7 ± 0.0	0.67 ± 0.01	0.77 ± 0.0	0.72 ± 0.01	1.19(3.48%)	1.4(-27.84%)	9.11(11.78%)	0.77 ± 0.06
		0.20	0.7 ± 0.0	0.67 ± 0.01	0.77 ± 0.0	0.72 ± 0.01	1.12(-2.61%)	1.3(-32.99%)	8.97(10.06%)	0.68 ± 0.04
	L1LR	0.05	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.14(-8.8%)	1.62(-29.87%)	5.67(11.39%)	0.29 ± 0.02
		0.10	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.14(-8.8%)	1.55(-32.9%)	5.85(14.93%)	0.29 ± 0.01
		0.20	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.14(-8.8%)	1.5(-35.06%)	5.87(15.32%)	0.28 ± 0.01
	L2LR	0.05	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.17(-7.14%)	1.92(-20.33%)	6.36(22.54%)	0.27 ± 0.01
		0.10	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.17(-7.14%)	1.85(-23.24%)	6.27(20.81%)	0.27 ± 0.01
		0.20	0.68 ± 0.0	0.67 ± 0.01	0.73 ± 0.0	0.72 ± 0.01	1.17(-6.35%)	1.68(-30.29%)	6.26(20.62%)	0.29 ± 0.01
	MLP	0.05	0.69 ± 0.01	0.67 ± 0.01	0.74 ± 0.01	0.72 ± 0.01	1.2(-11.11%)	1.67(-27.39%)	8.2(26.35%)	0.39 ± 0.07
		0.10	0.69 ± 0.01	0.67 ± 0.01	0.74 ± 0.01	0.72 ± 0.01	1.2(-11.11%)	1.65(-28.26%)	8.19(26.19%)	0.41 ± 0.06
		0.20	0.69 ± 0.01	0.67 ± 0.01	0.74 ± 0.01	0.72 ± 0.01	1.2(-10.37%)	1.62(-29.57%)	8.36(28.81%)	0.42 ± 0.07
Diabetes	GBDT	0.05	0.65 ± 0.0	0.64 ± 0.0	0.7 ± 0.0	0.7 ± 0.0	1.37(-3.6%)	1.16(-58.87%)	-4521.05(-92.44%)	50.03 ± 8.06
		0.10	0.65 ± 0.0	0.64 ± 0.0	0.7 ± 0.0	0.7 ± 0.0	1.36(-2.16%)	1.35(-52.13%)	-5505.82(-90.8%)	58.29 ± 7.65
		0.20	0.65 ± 0.0	0.64 ± 0.0	0.7 ± 0.0	0.7 ± 0.0	1.35(-2.88%)	1.46(-48.23%)	-5258.28(-91.21%)	54.67 ± 7.11
	L1LR	0.05	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.2(-25.93%)	2.31(-11.15%)	-11250.28(46.0%)	5.23 ± 0.68
		0.10	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.2(-25.93%)	2.31(-11.15%)	-11190.99(46.29%)	5.3 ± 0.7
		0.20	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.2(-25.93%)	2.31(-11.15%)	-7913.34(62.02%)	5.09 ± 0.63
	L2LR	0.05	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.24(-23.46%)	2.31(-14.44%)	-23047.62(22.58%)	7.05 ± 1.0
		0.10	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.24(-23.46%)	2.31(-14.44%)	-23047.64(22.58%)	7.12 ± 0.99
		0.20	0.62 ± 0.0	0.62 ± 0.0	0.66 ± 0.0	0.66 ± 0.0	1.24(-23.46%)	2.31(-14.44%)	-14691.43(21.86%)	7.41 ± 0.64
	MLP	0.05	0.65 ± 0.01	0.63 ± 0.0	0.71 ± 0.01	0.68 ± 0.0	1.41(-13.5%)	1.73(-35.45%)	-46675.04(33.81%)	40.41 ± 30.18
		0.10	0.65 ± 0.01	0.63 ± 0.0	0.71 ± 0.01	0.68 ± 0.0	1.41(-13.5%)	1.72(-35.82%)	-46689.47(33.84%)	38.03 ± 27.63
		0.20	0.65 ± 0.01	0.63 ± 0.0	0.71 ± 0.01	0.68 ± 0.0	1.39(-14.72%)	1.73(-35.45%)	-47723.79(4.23%)	30.72 ± 19.28
FICO	GBDT	0.05	0.77 ± 0.0	0.72 ± 0.01	0.85 ± 0.0	0.79 ± 0.01	1.97(-44.97%)	0.87(-51.93%)	-58.85(20.61%)	132.34 ± 34.38
		0.10	0.77 ± 0.0	0.72 ± 0.01	0.85 ± 0.0	0.79 ± 0.01	2.03(-43.3%)	0.89(-50.83%)	-58.47(21.13%)	162.91 ± 37.45
		0.20	0.77 ± 0.0	0.72 ± 0.01	0.85 ± 0.0	0.79 ± 0.01	2.03(-42.18%)	0.88(-51.38%)	-56.13(24.28%)	163.64 ± 45.55
	L1LR	0.05	0.71 ± 0.0	0.7 ± 0.0	0.78 ± 0.0	0.77 ± 0.01	1.84(-25.31%)	1.89(4.42%)	-77.6(4.56%)	29.88 ± 6.18
		0.10	0.71 ± 0.0	0.7 ± 0.0	0.78 ± 0.0	0.77 ± 0.01	1.86(-24.7%)	1.96(8.29%)	-78.18(3.85%)	34.15 ± 7.9
		0.20	0.71 ± 0.0	0.7 ± 0.0	0.78 ± 0.0	0.77 ± 0.01	1.86(-24.7%)	2.09(15.47%)	-79.92(-1.71%)	42.69 ± 9.43
	L2LR	0.05	0.72 ± 0.0	0.71 ± 0.01	0.78 ± 0.0	0.77 ± 0.01	2.3(-16.36%)	1.8(-6.74%)	-57.96(12.02%)	285.3 ± 96.59
		0.10	0.72 ± 0.0	0.71 ± 0.01	0.78 ± 0.0	0.77 ± 0.01	2.28(17.09%)	1.79(-7.25%)	-57.11(10.38%)	303.19 ± 98.72
		0.20	0.72 ± 0.0	0.71 ± 0.01	0.78 ± 0.0	0.77 ± 0.01	2.24(-18.55%)	1.91(-1.04%)	-57.22(10.59%)	303.85 ± 97.78
	MLP	0.05	0.74 ± 0.01	0.71 ± 0.01	0.81 ± 0.01	0.78 ± 0.01	2.17(-18.11%)	1.63(-10.93%)	-79.53(15.44%)	124.03 ± 50.02
		0.10	0.74 ± 0.01	0.71 ± 0.01	0.81 ± 0.01	0.78 ± 0.01	2.18(-17.74%)	1.66(-9.29%)	-77.83(12.98%)	135.6 ± 56.71
		0.20	0.74 ± 0.01	0.71 ± 0.01	0.81 ± 0.01	0.78 ± 0.01	2.18(-17.74%)	1.71(-6.56%)	-78.07(13.33%)	156.08 ± 70.95
German Credit	GBDT	0.05	0.96 ± 0.01	0.75 ± 0.02	0.99 ± 0.0	0.77 ± 0.03	1.21(-12.95%)	2.13(13.9%)	-31442.17(58.53%)	6.28 ± 3.44
		0.10	0.96 ± 0.01	0.75 ± 0.02	0.99 ± 0.0	0.77 ± 0.03	1.21(-12.95%)	1.8(-3.74%)	-31253.08(58.78%)	6.87 ± 3.83
		0.20	0.96 ± 0.01	0.75 ± 0.02	0.99 ± 0.0	0.77 ± 0.03	1.2(-12.23%)	1.91(2.14%)	-36087.77(52.4%)	7.78 ± 4.46
	L1LR	0.05	0.75 ± 0.01	0.75 ± 0.02	0.8 ± 0.01	0.78 ± 0.04	1.03(-20.77%)	2.03(-17.14%)	-24474.67(61.9%)	0.79 ± 0.39
		0.10	0.75 ± 0.01	0.75 ± 0.02	0.8 ± 0.01	0.77 ± 0.04	1.04(-20.0%)	2.01(-17.96%)	-24862.18(-61.3%)	0.79 ± 0.38
		0.20	0.75 ± 0.01	0.75 ± 0.02	0.8 ± 0.01	0.78 ± 0.04	1.03(-20.77%)	2.12(-13.47%)	-25849.27(-59.76%)	0.7 ± 0.17
	L2LR	0.05	0.78 ± 0.01	0.76 ± 0.03	0.83 ± 0.01	0.79 ± 0.04	1.17(-22.52%)	3.0(-2.6%)	-40660.55(63.68%)	2.05 ± 1.58
		0.10	0.78 ± 0.01	0.76 ± 0.03	0.83 ± 0.01	0.79 ± 0.04	1.18(-21.85%)	3.03(-1.62%)	-40228.76(64.06%)	1.84 ± 1.02
		0.20	0.78 ± 0.01	0.76 ± 0.03	0.83 ± 0.01	0.79 ± 0.04	1.17(-22.52%)	2.93(-4.87%)	-40136.71(64.15%)	1.71 ± 0.82
	MLP	0.05	0.81 ± 0.04	0.76 ± 0.03	0.87 ± 0.04	0.78 ± 0.04	1.25(-21.88%)	2.57(-4.46%)	-46257.34(61.31%)	2.99 ± 1.42
		0.10	0.81 ± 0.05	0.76 ± 0.03	0.87 ± 0.04	0.78 ± 0.04	1.23(-23.13%)	2.56(-4.83%)	-46884.11(60.79%)	3.04 ± 1.67
		0.20	0.81 ± 0.04	0.76 ± 0.03	0.87 ± 0.04	0.78 ± 0.04	1.21(-24.38%)	2.6(-3.35%)	-41223.18(65.52%)	2.55 ± 1.47
Headline	GBDT	0.05	0.82 ± 0.0	0.81 ± 0.0	0.9 ± 0.0	0.89 ± 0.0	1.74(-4.4%)	2.49(5.96%)	-407.77(-3.13%)	22.98 ± 8.46
		0.10	0.82 ± 0.0	0.81 ± 0.0	0.9 ± 0.0	0.89 ± 0.0	1.71(-6.04%)	2.51(6.81%)	-432.26(-9.32%)	20.88 ± 7.71
		0.20	0.82 ± 0.0	0.81 ± 0.0	0.9 ± 0.0	0.89 ± 0.0	1.53(-15.93%)	2.22(-5.53%)	-543.65(-37.49%)	8.83 ± 2.41
	L1LR	0.05	0.78 ± 0.0	0.78 ± 0.0	0.85 ± 0.0	0.85 ± 0.0	1.54(-19.79%)	2.94(17.13%)	-576.99(-3.25%)	3.97 ± 0.15
		0.10	0.78 ± 0.0	0.78 ± 0.0	0.85 ± 0.0	0.85 ± 0.0	1.55(-19.27%)	2.94(17.13%)	-577.03(-3.26%)	4.16 ± 0.17
		0.20	0.78 ± 0.0	0.78 ± 0.0	0.85 ± 0.0	0.85 ± 0.0	1.47(-23.44%)	2.94(17.13%)	-577.7(-3.38%)	2.54 ± 0.12
	L2LR	0.05	0.78 ± 0.0	0.78 ± 0.0	0.86 ± 0.0	0.85 ± 0.0	1.59(-19.7%)	2.94(17.6%)	-556.65(0.13%)	4.81 ± 0.2
		0.10	0.78 ± 0.0	0.78 ± 0.0	0.85 ± 0.0	0.85 ± 0.0	1.6(-19.19%)	2.94(17.6%)	-573.97(-3.24%)	5.1 ± 0.25
		0.20	0.78 ± 0.0	0.78 ± 0.0	0.85 ± 0.0	0.85 ± 0.0	1.5(-24.24%)	2.94(17.6%)	-574.67(-3.37%)	3.22 ± 0.13
	MLP	0.05	0.83 ± 0.01	0.81 ± 0.0	0.91 ± 0.01	0.89 ± 0.0	1.64(-19.21%)	1.97(-14.72%)	-617.43(-25.15%)	7.02 ± 1.86
		0.10	0.83 ± 0.01	0.81 ± 0.0	0.91 ± 0.01	0.89 ± 0.0	1.64(-19.21%)	1.97(-14.72%)	-604.44(-22.51%)	7.47 ± 2.23
		0.20	0.83 ± 0.01	0.81 ± 0.0	0.91 ± 0.01	0.89 ± 0.0	1.5(-26.11%)	2.06(-10.82%)	-570.13(-15.56%)	4.1 ± 0.79
MIMIC	GBDT	0.05	0.91 ± 0.0	0.9 ± 0.0	0.87 ± 0.0	0.85 ± 0.0	1.21(2.54%)	0.52(-59.38%)	-19.06(-0.74%)	2.93 ± 0.39
		0.10	0.91 ± 0.0	0.9 ± 0.0	0.87 ± 0.0	0.85 ± 0.0	1.21(2.54%)	0.48(-62.5%)	-19.08(-0.85%)	2.98 ± 0.39
		0.20	0.91 ± 0.0	0.9 ± 0.0	0.87 ± 0.0	0.85 ± 0.0	1.21(2.54%)	0.41(-67.97%)	-18.86(0.32%)	3.32 ± 0.43
	L1LR	0.05	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.17(1.74%)	1.11(-75.5%)	-21.32(-7.89%)	0.75 ± 0.06
		0.10	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.18(2.61%)	1.15(-74.61%)	-21.48(-8.7%)	0.77 ± 0.07
		0.20	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.18(2.61%)	1.15(-74.61%)	-21.48(-8.7%)	0.79 ± 0.08
	L2LR	0.05	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.19(2.59%)	1.15(-73.5%)	-21.37(-8.7%)	0.86 ± 0.1
		0.10	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.19(2.59%)	1.15(-73.5%)	-21.41(-8.9%)	0.84 ± 0.09
		0.20	0.89 ± 0.0	0.89 ± 0.0	0.8 ± 0.0	0.8 ± 0.0	1.19(2.59%)	1.15(-73.5%)	-21.48(-9.26%)	0.91 ± 0.09
	MLP	0.05	0.9 ± 0.0	0.9 ± 0.0	0.87 ± 0.01	0.85 ± 0.0	1.21(2.54%)	0.58(-72.12%)	-18.22(-5.62%)	1.35 ± 0.15
		0.10	0.9 ± 0.0	0.9 ± 0.0	0.87 ± 0.01	0.8				

748 I All-Opt⁻ Variants Performance

749 In this section, we will mainly show the model performance of All-Opt[©] and All-Opt¹, which are the
750 two gradient-based optimization methods used for SEV[©] and SEV¹ optimization. Table 10 shows the
751 SEV¹, ℓ_∞ and model performance after applying All-Opt¹ methods for different models on different
752 datasets with different levels of flexibility. It is evident that All-Opt^F has provided a significant
753 decrease in SEV, so that its values are close to 1, providing much sparser explanations without model
754 performance loss and closeness/credibility loss in explanations. Similar findings are observed in
755 Table 11.

Table 10: The model performance for All-Opt¹

DATASET	MODEL	TRAIN ACCURACY	TEST ACCURACY	TRAIN AUC	TEST AUC	MEAN SEV ⁻	MEAN ℓ_∞	TRAINING TIME(S)	MEAN LOG-LIKELIHOOD
Adult	GBDT	0.87 ± 0.02	0.84 ± 0.02	0.93 ± 0.01	0.90 ± 0.01	1.00 ± 0.00	5.67 ± 0.34	2010 ± 24	-39654.89 ± 4201.17
	LR	0.84 ± 0.01	0.84 ± 0.01	0.90 ± 0.02	0.89 ± 0.01	1.03 ± 0.01	3.21 ± 0.02	60 ± 1	-70566.06 ± 10678.32
	MLP	0.86 ± 0.01	0.85 ± 0.01	0.91 ± 0.02	0.91 ± 0.01	1.00 ± 0.00	9.52 ± 1.45	82 ± 3	-58049.77 ± 9932.16
COMPAS	GBDT	0.70 ± 0.01	0.68 ± 0.01	0.74 ± 0.01	0.71 ± 0.01	1.01 ± 0.01	1.50 ± 0.04	244 ± 4	10.74 ± 0.98
	LR	0.68 ± 0.01	0.68 ± 0.02	0.74 ± 0.01	0.73 ± 0.02	1.00 ± 0.01	2.13 ± 0.01	11 ± 1	9.17 ± 1.02
	MLP	0.68 ± 0.01	0.67 ± 0.02	0.74 ± 0.02	0.72 ± 0.01	1.01 ± 0.01	1.90 ± 0.11	16 ± 1	14.57 ± 1.23
Diabetes	GBDT	0.62 ± 0.01	0.63 ± 0.01	0.62 ± 0.01	0.64 ± 0.01	1.07 ± 0.01	1.78 ± 0.34	10548 ± 324	-14013.49 ± 2784.36
	LR	0.62 ± 0.04	0.62 ± 0.04	0.63 ± 0.01	0.63 ± 0.01	1.07 ± 0.00	1.39 ± 0.01	217 ± 3	-40190.09 ± 10453.69
	MLP	0.62 ± 0.01	0.65 ± 0.01	0.65 ± 0.01	0.64 ± 0.02	1.07 ± 0.00	2.50 ± 0.32	318 ± 5	-18013.49 ± 3894.36
FICO	GBDT	0.70 ± 0.02	0.70 ± 0.02	0.77 ± 0.01	0.77 ± 0.02	1.19 ± 0.10	0.84 ± 0.12	864 ± 23	-40.44 ± 4.32
	LR	0.70 ± 0.02	0.70 ± 0.02	0.77 ± 0.01	0.77 ± 0.02	1.10 ± 0.10	1.91 ± 0.33	19 ± 1	-20.32 ± 0.18
	MLP	0.72 ± 0.01	0.72 ± 0.01	0.78 ± 0.02	0.78 ± 0.01	1.28 ± 0.09	1.23 ± 0.21	28 ± 0	-26.04 ± 0.43
German Credit	GBDT	0.94 ± 0.02	0.73 ± 0.02	0.99 ± 0.01	0.76 ± 0.02	1.02 ± 0.01	1.21 ± 0.05	99 ± 1	-27701.04 ± 3431.99
	LR	0.77 ± 0.01	0.75 ± 0.01	0.82 ± 0.02	0.77 ± 0.01	1.00 ± 0.00	1.39 ± 0.05	2 ± 0	-58065.80 ± 6843.21
	MLP	0.82 ± 0.01	0.73 ± 0.03	0.93 ± 0.02	0.75 ± 0.02	1.00 ± 0.00	1.17 ± 0.08	3 ± 1	-85816.95 ± 13728.23
Headline	GBDT	0.80 ± 0.01	0.76 ± 0.02	0.90 ± 0.01	0.89 ± 0.01	1.04 ± 0.02	2.45 ± 0.57	2732 ± 101	-4.37 ± 1.28
	LR	0.77 ± 0.01	0.78 ± 0.01	0.86 ± 0.01	0.85 ± 0.01	1.00 ± 0.01	2.77 ± 0.44	78 ± 0	-2.39 ± 0.11
	MLP	0.76 ± 0.02	0.77 ± 0.03	0.87 ± 0.02	0.86 ± 0.02	1.03 ± 0.03	2.78 ± 0.13	102 ± 1	-2.57 ± 0.89
MIMIC	GBDT	0.88 ± 0.01	0.88 ± 0.01	0.84 ± 0.01	0.82 ± 0.02	1.06 ± 0.04	3.66 ± 0.02	2799 ± 102	-16.36 ± 0.54
	LR	0.88 ± 0.01	0.88 ± 0.01	0.84 ± 0.01	0.82 ± 0.02	1.03 ± 0.03	3.67 ± 0.72	87 ± 2	-17.77 ± 2.22
	MLP	0.89 ± 0.01	0.89 ± 0.02	0.84 ± 0.03	0.82 ± 0.03	1.00 ± 0.00	1.29 ± 0.20	115 ± 2	-10.38 ± 3.87

Table 11: The model performance for All-Opt[©]

DATASET	MODEL	TRAIN ACCURACY	TEST ACCURACY	TRAIN AUC	TEST AUC	MEAN SEV [©]	MEAN ℓ_∞	MEAN LOG-LIKELIHOOD
Adult	GBDT	0.90 ± 0.00	0.83 ± 0.01	0.89 ± 0.01	0.89 ± 0.01	1.14 ± 0.03	1.87 ± 0.03	289.07 ± 52.79
	LR	0.84 ± 0.00	0.84 ± 0.01	0.91 ± 0.01	0.90 ± 0.01	1.01 ± 0.01	2.56 ± 0.43	299.04 ± 17.24
	MLP	0.85 ± 0.01	0.84 ± 0.01	0.92 ± 0.01	0.91 ± 0.01	1.00 ± 0.02	2.37 ± 0.19	297.14 ± 32.16
COMPAS	GBDT	0.68 ± 0.01	0.68 ± 0.01	0.72 ± 0.01	0.74 ± 0.02	1.02 ± 0.02	1.34 ± 0.47	10.28 ± 2.14
	LR	0.68 ± 0.01	0.68 ± 0.01	0.72 ± 0.01	0.74 ± 0.02	1.00 ± 0.00	2.49 ± 0.21	8.67 ± 1.32
	MLP	0.67 ± 0.01	0.67 ± 0.02	0.74 ± 0.01	0.72 ± 0.01	1.05 ± 0.05	1.92 ± 0.05	7.22 ± 0.56
Diabetes	GBDT	0.62 ± 0.01	0.62 ± 0.02	0.66 ± 0.01	0.66 ± 0.02	1.05 ± 0.00	1.99 ± 0.01	-5231.53 ± 489.52
	LR	0.62 ± 0.01	0.62 ± 0.02	0.66 ± 0.01	0.66 ± 0.02	1.05 ± 0.00	2.89 ± 0.46	-5937.66 ± 638.77
	MLP	0.62 ± 0.01	0.62 ± 0.01	0.67 ± 0.01	0.67 ± 0.01	1.05 ± 0.00	2.12 ± 0.01	-5217.39 ± 497.78
FICO	GBDT	0.70 ± 0.01	0.70 ± 0.00	0.78 ± 0.01	0.78 ± 0.01	1.48 ± 0.09	0.90 ± 0.01	-55.09 ± 6.79
	LR	0.70 ± 0.01	0.70 ± 0.00	0.78 ± 0.01	0.78 ± 0.01	1.41 ± 0.08	1.60 ± 0.27	-15.66 ± 7.01
	MLP	0.70 ± 0.01	0.69 ± 0.11	0.79 ± 0.02	0.78 ± 0.02	1.28 ± 0.19	1.23 ± 0.05	-18.47 ± 8.98
German Credit	GBDT	0.75 ± 0.01	0.76 ± 0.01	0.82 ± 0.01	0.80 ± 0.01	1.00 ± 0.00	1.00 ± 0.00	-15797.31 ± 2134.01
	LR	0.75 ± 0.01	0.76 ± 0.01	0.82 ± 0.01	0.80 ± 0.01	1.00 ± 0.00	1.00 ± 0.00	-45070.76 ± 7924.23
	MLP	0.86 ± 0.02	0.79 ± 0.01	0.92 ± 0.01	0.80 ± 0.01	1.00 ± 0.00	1.00 ± 0.00	-30917.95 ± 5534.23
Headline	GBDT	0.78 ± 0.02	0.79 ± 0.01	0.85 ± 0.01	0.85 ± 0.01	1.26 ± 0.03	-1.72 ± 0.01	-4.20 ± 2.97
	LR	0.78 ± 0.02	0.79 ± 0.01	0.85 ± 0.01	0.85 ± 0.01	1.29 ± 0.10	2.93 ± 0.02	-2.93 ± 1.28
	MLP	0.78 ± 0.02	0.78 ± 0.03	0.84 ± 0.01	0.84 ± 0.01	1.15 ± 0.12	1.69 ± 0.16	-2.87 ± 1.51
MIMIC	GBDT	0.90 ± 0.01	0.89 ± 0.01	0.80 ± 0.00	0.80 ± 0.00	1.05 ± 0.05	1.00 ± 0.00	-21.80 ± 2.45
	LR	0.90 ± 0.01	0.89 ± 0.01	0.80 ± 0.00	0.80 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	-28.74 ± 0.75
	MLP	0.89 ± 0.01	0.89 ± 0.01	0.84 ± 0.01	0.81 ± 0.00	1.01 ± 0.01	0.06 ± 0.01	-29.35 ± 0.36

756 J SEV^T in tree-based models

757 In this section, we show the model performance and SEV^T values for different types of tree-based
 758 models. As discussed in section 4.2, the similarity and closeness metrics in SEV^T are all ℓ_0 norm, so
 759 we only need to compute the mean SEV^T for each tree. Table 12 shows that most of the tree-based
 760 models can provide sparse explanations ($SEV^T \leq 2$), and we can also find a decision tree with the
 761 same model performance as the other tree-based models from $SEV^T=1$ to TOpt.

Table 12: The model performance with different tree-based methods

DATASET	METHODS	TRAIN ACC	TEST ACC	MEAN SEV^T
Adult	CART	0.84 ± 0.01	0.84 ± 0.01	1.11 ± 0.01
	C4.5	0.85 ± 0.01	0.84 ± 0.00	1.10 ± 0.02
	GOSDT	0.81 ± 0.01	0.81 ± 0.01	1.08 ± 0.01
	Topt	0.82 ± 0.01	0.82 ± 0.01	1.00 ± 0.00
COMPAS	CART	0.68 ± 0.00	0.65 ± 0.01	1.02 ± 0.01
	C4.5	0.68 ± 0.00	0.65 ± 0.01	1.02 ± 0.01
	GOSDT	0.67 ± 0.02	0.65 ± 0.01	1.12 ± 0.02
	Topt	0.66 ± 0.01	0.67 ± 0.01	1.00 ± 0.00
Diabetes	CART	0.63 ± 0.01	0.63 ± 0.01	1.00 ± 0.00
	C4.5	0.63 ± 0.01	0.63 ± 0.01	1.00 ± 0.00
	GOSDT	0.61 ± 0.01	0.60 ± 0.01	1.00 ± 0.00
	Topt	0.62 ± 0.01	0.63 ± 0.01	1.00 ± 0.00
FICO	CART	0.71 ± 0.01	0.71 ± 0.01	1.10 ± 0.03
	C4.5	0.71 ± 0.01	0.71 ± 0.01	1.13 ± 0.05
	GOSDT	0.70 ± 0.01	0.69 ± 0.01	1.80 ± 0.02
	Topt	0.70 ± 0.01	0.71 ± 0.01	1.00 ± 0.02
German Credit	CART	0.75 ± 0.01	0.70 ± 0.01	1.00 ± 0.02
	C4.5	0.75 ± 0.01	0.70 ± 0.01	1.00 ± 0.02
	GOSDT	0.75 ± 0.01	0.70 ± 0.01	1.00 ± 0.02
	Topt	0.75 ± 0.01	0.70 ± 0.01	1.00 ± 0.02
Headline	CART	0.78 ± 0.01	0.78 ± 0.00	1.27 ± 0.01
	C4.5	0.77 ± 0.01	0.77 ± 0.00	1.16 ± 0.02
	GOSDT	0.76 ± 0.01	0.76 ± 0.02	1.09 ± 0.02
	Topt	0.77 ± 0.00	0.77 ± 0.00	1.00 ± 0.00
MIMIC	CART	0.89 ± 0.01	0.89 ± 0.01	1.00 ± 0.00
	C4.5	0.89 ± 0.01	0.89 ± 0.01	1.00 ± 0.00
	GOSDT	0.89 ± 0.01	0.89 ± 0.01	1.00 ± 0.00
	Topt	0.89 ± 0.01	0.89 ± 0.01	1.00 ± 0.00

K The SEV¹ results after ExpO Optimization

For the ExpO comparison experiment, we used the fidelity metrics from Plumb et al. [2020] as the penalty term for regularizing the original model. Then we evaluated the optimized model with SEV⁻. We used two kinds of fidelity metrics as the regularization term: 1D fidelity and 1D fidelity. Both of these two penalty terms aim to optimize the model f such that the local model g [Ribeiro et al., 2016b, Plumb et al., 2018] accurately approximates f in the neighborhood N_x , which is equivalent to minimizing:

$$\ell_{\text{fed}}(f, g, N_x) = \mathbb{E}_{\mathbf{x}' \sim N_x} [g(\mathbf{x}') - f(\mathbf{x}')]^2. \quad (10)$$

The local model g 's are linear models, and the N_x are points sampled normally around the original query. The 1D version of Fidelity regularization requires sampling the points around each feature of x at a time, which saves time and computational complexity. Based on the above equation, we rewrite the overall objective function as:

$$\min_{f \in \mathcal{F}} \ell_{\text{BCE}} + C_F \ell_{\text{fed}} \quad (11)$$

where ℓ_{BCE} is the Binary Cross Entropy Loss to control the accuracy of the training model, C_F is the strength of the fidelity term, and the training process is the same All-Opt⁻ optimization, which we used 80 epochs for basic training process, 20 epochs for regularization.

In this section, we show the SEV⁻ and training time for ExpO regularizer in **LR** and **MLP** models with 1D Fidelity (1DFed) and Global Fidelity (Fed) regularizers. Comparing the mean SEV¹ of Table 13 with Table 7, it is evident that with the optimization through Fed or 1DFed, the optimized models do not provide sparse explanations. In addition, it takes a long time to calculate Fed and 1DFed since the regularizer's complexity is determined by the number of queries, features, as well as the points samples around the queries. For SEV⁻, the complexity is determined only by the number of queries and the number of features, so it is much easier to calculate.

Table 13: Model performance, SEV¹ and training time of LR and MLPs after ExpO with different datasets

DATASET	MODEL	REGULARIZER	TRAIN ACCURACY	TEST ACCURACY	TRAIN AUC	TEST AUC	MEAN SEV ¹	TRAINING TIME(s)
Adult	LR	Fed	0.85 ± 0.01	0.84 ± 0.01	0.90 ± 0.01	0.89 ± 0.01	1.23 ± 0.02	1350 ± 162
	LR	1DFed	0.84 ± 0.02	0.84 ± 0.01	0.90 ± 0.01	0.90 ± 0.02	1.17 ± 0.02	510 ± 23
	MLP	Fed	0.85 ± 0.01	0.83 ± 0.02	0.90 ± 0.01	0.89 ± 0.01	1.27 ± 0.02	1580 ± 50
	MLP	1DFed	0.85 ± 0.01	0.83 ± 0.02	0.90 ± 0.01	0.89 ± 0.01	1.27 ± 0.02	686 ± 23
COMPAS	LR	Fed	0.67 ± 0.02	0.66 ± 0.01	0.72 ± 0.02	0.72 ± 0.02	1.22 ± 0.04	58 ± 10
	LR	1DFed	0.65 ± 0.02	0.65 ± 0.01	0.73 ± 0.01	0.72 ± 0.02	1.27 ± 0.02	90 ± 5
	MLP	Fed	0.68 ± 0.02	0.66 ± 0.01	0.74 ± 0.02	0.72 ± 0.01	1.28 ± 0.03	125 ± 14
	MLP	1DFed	0.66 ± 0.02	0.66 ± 0.02	0.72 ± 0.02	0.71 ± 0.01	1.28 ± 0.2	128 ± 15
Diabetes	LR	Fed	0.63 ± 0.02	0.62 ± 0.01	0.60 ± 0.02	0.60 ± 0.01	1.50 ± 0.01	3625 ± 412
	LR	1DFed	0.63 ± 0.02	0.62 ± 0.01	0.60 ± 0.02	0.60 ± 0.01	1.46 ± 0.01	1842 ± 245
	MLP	Fed	0.63 ± 0.02	0.62 ± 0.01	0.60 ± 0.02	0.60 ± 0.01	1.52 ± 0.01	4372 ± 316
	MLP	1DFed	0.63 ± 0.02	0.62 ± 0.01	0.60 ± 0.02	0.60 ± 0.01	1.46 ± 0.01	2032 ± 124
FICO	LR	Fed	0.71 ± 0.01	0.71 ± 0.01	0.78 ± 0.02	0.78 ± 0.01	2.76 ± 0.12	150 ± 21
	LR	1DFed	0.71 ± 0.02	0.71 ± 0.01	0.77 ± 0.01	0.78 ± 0.01	2.76 ± 0.21	150 ± 14
	MLP	Fed	0.72 ± 0.02	0.71 ± 0.01	0.79 ± 0.02	0.78 ± 0.02	2.67 ± 0.14	210 ± 13
	MLP	1DFed	0.72 ± 0.02	0.71 ± 0.01	0.78 ± 0.02	0.77 ± 0.02	2.80 ± 0.35	195 ± 14
German Credit	LR	Fed	0.78 ± 0.02	0.76 ± 0.01	0.82 ± 0.02	0.80 ± 0.01	1.65 ± 0.12	28 ± 0
	LR	1DFed	0.77 ± 0.02	0.73 ± 0.02	0.80 ± 0.01	0.76 ± 0.02	1.76 ± 0.02	15 ± 0
	MLP	Fed	0.75 ± 0.02	0.72 ± 0.02	0.82 ± 0.01	0.78 ± 0.02	1.70 ± 0.03	33 ± 2
	MLP	1DFed	0.70 ± 0.00	0.70 ± 0.00	0.72 ± 0.02	0.73 ± 0.01	1.70 ± 0.03	20 ± 0
Headline	LR	Fed	0.77 ± 0.04	0.77 ± 0.01	0.85 ± 0.01	0.85 ± 0.00	1.87 ± 0.01	680 ± 21
	LR	1DFed	0.77 ± 0.01	0.77 ± 0.01	0.84 ± 0.01	0.85 ± 0.01	1.87 ± 0.02	562 ± 32
	MLP	Fed	0.77 ± 0.02	0.78 ± 0.01	0.85 ± 0.02	0.85 ± 0.03	1.87 ± 0.04	762 ± 56
	MLP	1DFed	0.77 ± 0.02	0.77 ± 0.01	0.84 ± 0.02	0.85 ± 0.01	1.87 ± 0.04	852 ± 72
MIMIC	LR	Fed	0.89 ± 0.02	0.89 ± 0.02	0.77 ± 0.01	0.77 ± 0.01	1.18 ± 0.02	712 ± 42
	LR	1DFed	0.89 ± 0.02	0.88 ± 0.01	0.78 ± 0.02	0.77 ± 0.02	1.17 ± 0.02	646 ± 42
	MLP	Fed	0.88 ± 0.00	0.88 ± 0.00	0.78 ± 0.00	0.77 ± 0.01	1.15 ± 0.01	960 ± 27
	MLP	1DFed	0.88 ± 0.01	0.88 ± 0.01	0.78 ± 0.01	0.78 ± 0.01	1.16 ± 0.01	873 ± 18

783 L Proof of Theorem 4.1

784 **Theorem L.1.** *With a single decision classifier DT and a positively-predicted query x_i , define N_i*
 785 *as the leaf that captures it. If N_i has a sibling leaf, or any internal node in its decision path has a*
 786 *negatively-predicted child leaf, then SEV^T is equal to 1.*

787 SEV^- is defined as the number of features that need to change within the given classification tree. If
 788 you have switched a particular node from one path to another, it adds one to SEV^- . Therefore, for
 789 the internal nodes along the SEV^- path, if N_i has a sibling leaf node, if we goes up to its parent node
 790 and goes the opposite direction to change the query value for counterfactual explanation, the modified
 791 instance will be directly predicted as negative, which leads to SEV^- being equal to 1 in this case.

792 Figure 11 shows an example for SEV^T being exactly 1, and a case illustrating that if N does not have
 793 a sibling or any internal node in its decision path that has a negatively-predicted child leaf, SEV^T
 794 should be greater than or equal to 1. In Figure 11, the left trees are the full decision trees, where the
 795 blue nodes are the negatively predicted leaf nodes and the red ones are positively predicted. The red
 796 arrows graph represents the decision path for a specific instance. The person icon with a plus sign is
 797 N_i that we would like to calculate SEV^T on. The right tree is the subtree of the left tree. The person
 798 icon with a minus is the query and the blue arrows indicate a decision pathway for SEV Explanation.

799 If the query is predicted as positive in node ④, it is easy to see that if we go up to node ③ and goes
 800 the opposite direction as the decision path for x_i , then you can directly get a negative prediction.
 801 In other words, if you change the feature C in the query to make it doesn't satisfy the node ③'s
 802 condition, then it can be prediction as negative, which means that $SEV^T=1$.

803 For $SEV^T \geq 1$ case, if the query predicted as positive in node ⑦, since it does not have a sibling
 804 leaf node, then if it goes to its parent node ⑤ and goes the opposite direction, then it would reach
 805 node ⑥. However, if we don't know the query x_i 's value, then I am unable to know whether I need
 806 to change the condition in node ⑥ for higher SEV^T . Therefore, in this case SEV^T can be only
 807 guaranteed to be greater or equal to 1.

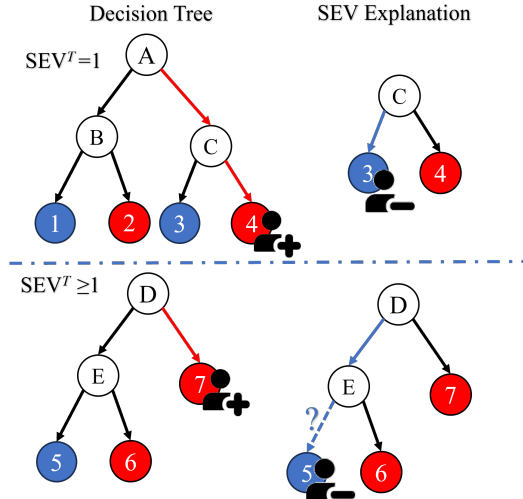


Figure 11: Example of $SEV^T=1$ in Theorem 4.1

808 M Proof of Theorem 4.2

809 **Theorem M.1.** *With a single decision tree classifier DT and a positively-predicted query x_i , with*
 810 *the set of all negatively predicted leaves as reference points, both SEV^- and the ℓ_0 distance (edit*
 811 *distance) between the query and the SEV^- explanation is minimized.*

812 Proof (Optimality of Explanation Path):

813 The definition for SEV^- is the minimum number of features that is needed for a positively predicted
 814 query x_i to aligned with the reference point in order to be predicted as negative. For tree-based
 815 classifiers, the decisions are all made in the leaf nodes. Since we have set of all the negatively
 816 predicted leaves as the reference points, then the ℓ_0 distance (edit distance) between the query and the
 817 SEV^- explanation is equivalent to be the minimum ℓ_0 distance between the query and the negatively
 818 predicted leaf nodes. Each node can be considered as a list of rules of conditions that needs to be
 819 satisfied. If a query would like to be predicted as negative in a specific node, then it needs to change
 820 some of the feature values in the query so as to be predicted as negative, and the number of changed
 821 feature is SEV^- . Therefore, SEV^- and the ℓ_0 distance are the same in this theorem.

822 Next, we would like to show that if one of the negatively predicted leaf nodes is not considered
 823 as reference point, then SEV^- is not minimized. It is really easy to give an counterexample: if
 824 we have a decision tree shown in Figure 12 with white nodes as root/internal nodes, blue nodes
 825 as negatively predicted node, and the red ones as positively predicted. Suppose we have a query
 826 predicted as positive, with feature values $\{A : \text{False}, B : \text{False}, C : \text{False}\}$, and only regard node ①
 827 as the reference point, then both feature A and C should be change to True, in order to do a negative
 828 prediction, in other words, if only node ① is the reference point, then $SEV^- = 2$. However, based on
 829 Theorem 4.1, since node ④ has a sibling leaf predicted as negative, then the SEV^- is not minimized.

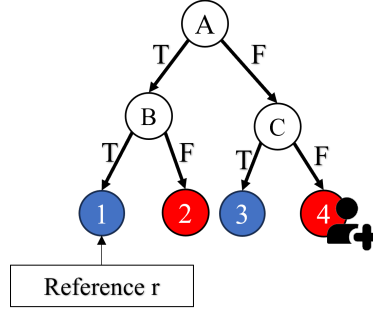


Figure 12: An counterexample with fewer reference point

830 Lastly, we would like to show that with all the negative leaf nodes considered as reference points if
 831 an new reference points is added, the SEV^- cannot be further minimized. Since we know that the
 832 reference points should be predicted as negative, so the newly aded reference should still belongs to
 833 one of the existing negative predicted leaf node, so SEV^- cannot be further minimized.

834 To sum up, we have proved that with the set of all negatively predicted leaves as reference points, both
 835 SEV^- and the ℓ_0 distance (edit distance) between the query and the SEV explanation is minimized.

Table 14: Different SEV Variants Explanations in MIMIC datasets

	PREICULOS	GCS	HEARTRATE_MAX	MEANBP_MIN	RESPRATE_MIN	TEMPC_MIN	URINEOUTPUT
Query	43806.28	10.00	91.00	29.00	9.00	34.50	162.98
SEV-I	2215.88	—	—	—	—	—	—
SEV-F	2215.88	—	—	—	—	—	—
SEV-C	8739.30	—	—	—	—	—	—
SEV-T	—	—	—	—	—	—	595.48
Query	0.51	15.00	105.00	21.00	20.00	32.28	7.98
SEV-I	—	—	—	59.35	—	—	—
SEV-F	—	—	—	59.35	—	—	—
SEV-C	—	—	—	56.95	—	36.11	—
SEV-T	—	—	—	—	—	—	595.48
Query	1.34	3.00	139.00	33.00	11.00	35.56	247.98
SEV-I	—	13.89	—	—	—	—	—
SEV-F	—	13.89	—	—	—	—	—
SEV-C	—	9.24	105.96	59.24	—	—	—
SEV-T	—	—	—	—	—	—	595.48
Query	1.64	11.00	199.00	14.00	22.00	37.06	387.98
SEV-I	—	—	102.57	—	—	—	—
SEV-F	—	—	102.57	—	—	—	—
SEV-C	—	—	107.58	—	—	—	—
SEV-T	—	—	—	—	—	—	595.48
Query	6621.40	13.00	134.00	28.00	28.00	34.72	4.98
SEV-I	—	—	102.57	—	12.22	—	—
SEV-F	—	—	102.57	—	12.22	—	—
SEV-C	—	—	97.70	—	12.68	—	—
SEV-T	—	—	—	—	—	—	595.48

Table 15: Different SEV Variants Explanations in COMPAS datasets

	AGE	JUV_FEL_COUNT	JUV_MISD_COUNT	JUVENILE_CRIMES	PRIORS_COUNT
Query	50.00	0.00	0.00	0.00	11.00
SEV-I	—	—	—	—	2.21
SEV-F	—	—	—	—	2.21
SEV-C	—	—	—	—	4.63
SEV-T	—	—	—	—	2.50
Query	23.00	1.00	0.00	1.00	5.00
SEV-I	36.71	—	—	—	2.21
SEV-F	36.71	—	—	—	2.21
SEV-C	26.69	0.11	0.18	0.54	2.13
SEV-T	—	—	—	—	2.50
Query	21.00	0.00	2.00	3.00	3.00
SEV-I	—	—	—	0.12	—
SEV-F	—	—	—	0.12	—
SEV-C	26.69	—	—	0.54	—
SEV-T	33.50	—	—	—	—
Query	23.00	0.00	1.00	1.00	4.00
SEV-I	36.71	—	—	—	—
SEV-F	36.71	—	—	—	—
SEV-C	26.69	—	—	—	2.13
SEV-T	23.00	—	—	—	2.50
Query	21.00	0.00	0.00	0.00	1.00
SEV-I	36.71	—	—	—	—
SEV-F	36.71	—	—	—	—
SEV-C	28.02	—	—	—	—
SEV-T	22.50	—	—	—	—

Table 16: Different SEV Variants Explanations in FICO datasets

	EXTERNAL RISKESTIMATE	MSINCE OLDEST TRADEOPEN	MSINCE MOSTRECENT TRADEOPEN	AVERAGE MINFILE	SATISFACTORY TRADES	NUM TRADES	NUMTRADES 60EVER2	DEROGPUBREC	NUMTRADES90 EVER2	MAXDELTQ2 PUBLICREC LAST12M	NUMINQ LAST6M	NUMINQ LAST6 MEXCL7DAYS	NETFRACTION REVOLVING BURDEN
Query	60.00	Missing	8.00	88.00	55.00	0.00	0.00	0.00	0.00	4.00	1.00	1.00	54.00
SEV-I	72.21	---	---	---	---	---	---	---	---	---	---	---	---
SEV-F	72.21	---	---	---	---	---	---	---	---	---	---	---	---
SEV-C	70.82	---	---	---	---	---	---	---	---	---	---	---	---
SEV-T	74.50	---	---	---	---	---	---	---	---	---	---	---	---
Query	60.00	150.99	32.00	79.00	8.00	2.00	0.00	0.00	0.00	3.00	0.00	0.00	112.01
SEV-I	72.21	---	9.20	---	21.10	---	---	---	---	---	---	---	22.26
SEV-F	---	---	---	---	---	Missing	---	---	---	---	---	---	9.00
SEV-C	---	---	11.80	---	---	Missing	---	---	---	---	---	---	8.85
SEV-T	74.50	---	---	---	---	---	---	---	---	---	---	---	---
Query	60.00	197.00	17.00	81.00	16.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	6.00
SEV-I	72.21	---	---	---	---	---	---	---	---	---	---	---	---
SEV-F	72.21	---	---	---	---	---	---	---	---	---	---	---	---
SEV-C	---	---	---	---	---	Missing	---	---	---	---	---	---	---
SEV-T	74.50	---	---	---	---	---	---	---	---	---	---	---	---
Query	59.00	125.99	12.00	58.00	18.00	2.00	1.00	1.00	10.00	2.00	10.00	10.00	95.01
SEV-I	72.21	---	---	82.32	---	0.00	---	---	0.60	5.36	0.56	0.56	22.26
SEV-F	---	---	---	---	---	Missing	Missing	Missing	---	---	---	---	9.00
SEV-C	70.82	218.29	8.60	85.80	23.67	0.82	0.51	0.51	1.22	5.10	1.18	1.18	30.36
SEV-T	74.50	---	---	---	---	---	---	---	---	---	---	---	---
Query	69.00	280.01	11.00	125.00	16.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	45.00
SEV-I	---	---	---	---	---	---	---	---	---	5.36	---	---	---
SEV-F	---	---	---	---	---	---	---	---	---	5.36	---	---	---
SEV-C	---	---	---	---	---	---	---	---	---	5.10	---	---	---
SEV-T	74.50	---	---	---	---	---	---	---	---	---	---	---	---

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and precede the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Our motivation and claims are made within the abstract. We have provided experimental and theoretical results for cluster-based SEV, and its variants, and propose algorithm for improving the decision sparsity.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Yes, we have discuss the limitation of the work in the conclusion section.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.

- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: We have provided the theorem mostly for the tree-based SEV in the Section 4.2, and the corresponding proofs are shown in Appendix L and Appendix M.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Yes, all the experiment details are mentioned in the Appendix F. The detailed training process for the comparison with ExpO is shown in Appendix K.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.

- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: Yes, we have provided the code for training, and evaluation in the Experiment folder, and the script for running in Script folder.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

994 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
 995 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
 996 results?

997 Answer: [\[Yes\]](#)

998 Justification: Yes, we have already mentioned them in the Section F.

999 Guidelines:

- 1000 • The answer NA means that the paper does not include experiments.
- 1001 • The experimental setting should be presented in the core of the paper to a level of detail
- 1002 that is necessary to appreciate the results and make sense of them.
- 1003 • The full details can be provided either with the code, in appendix, or as supplemental
- 1004 material.

1005 **7. Experiment Statistical Significance**

1006 Question: Does the paper report error bars suitably and correctly defined or other appropriate
 1007 information about the statistical significance of the experiments?

1008 Answer: [\[Yes\]](#)

1009 Justification: Yes, all the training data has been run for 10 times, which is mentioned in
 1010 Section F, and all the results are calculated for error bars.

1011 Guidelines:

- 1012 • The answer NA means that the paper does not include experiments.
- 1013 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 1014 dence intervals, or statistical significance tests, at least for the experiments that support
- 1015 the main claims of the paper.
- 1016 • The factors of variability that the error bars are capturing should be clearly stated (for
- 1017 example, train/test split, initialization, random drawing of some parameter, or overall
- 1018 run with given experimental conditions).
- 1019 • The method for calculating the error bars should be explained (closed form formula,
- 1020 call to a library function, bootstrap, etc.)
- 1021 • The assumptions made should be given (e.g., Normally distributed errors).
- 1022 • It should be clear whether the error bar is the standard deviation or the standard error
- 1023 of the mean.
- 1024 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 1025 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 1026 of Normality of errors is not verified.
- 1027 • For asymmetric distributions, the authors should be careful not to show in tables or
- 1028 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 1029 error rates).
- 1030 • If error bars are reported in tables or plots, The authors should explain in the text how
- 1031 they were calculated and reference the corresponding figures or tables in the text.

1032 **8. Experiments Compute Resources**

1033 Question: For each experiment, does the paper provide sufficient information on the com-
 1034 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 1035 the experiments?

1036 Answer: [\[Yes\]](#)

1037 Justification: Yes, we have error bars for the time execution for each methods and the GPU
 1038 and CPU details in Appendix F.

1039 Guidelines:

- 1040 • The answer NA means that the paper does not include experiments.
- 1041 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 1042 or cloud provider, including relevant memory and storage.
- 1043 • The paper should provide the amount of compute required for each of the individual
- 1044 experimental runs as well as estimate the total compute.

- 1045 • The paper should disclose whether the full research project required more compute
1046 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1047 didn't make it into the paper).

1048 9. Code Of Ethics

1049 Question: Does the research conducted in the paper conform, in every respect, with the
1050 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1051 Answer: [Yes]

1052 Justification: Yes, the paper conforms, in every respect, with the NeurIPS Code of Ethics
1053 <https://neurips.cc/public/EthicsGuidelines>

1054 Guidelines:

- 1055 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1056 • If the authors answer No, they should explain the special circumstances that require a
1057 deviation from the Code of Ethics.
- 1058 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1059 eration due to laws or regulations in their jurisdiction).

1060 10. Broader Impacts

1061 Question: Does the paper discuss both potential positive societal impacts and negative
1062 societal impacts of the work performed?

1063 Answer: [Yes]

1064 Justification: Yes, we have mentioned the social impact in the conclusion. Our method has
1065 impact in that it provides sparser explanations for those subjected to decisions made by
1066 models, including in finance and criminal justice.

1067 Guidelines:

- 1068 • The answer NA means that there is no societal impact of the work performed.
- 1069 • If the authors answer NA or No, they should explain why their work has no societal
1070 impact or why the paper does not address societal impact.
- 1071 • Examples of negative societal impacts include potential malicious or unintended uses
1072 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1073 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1074 groups), privacy considerations, and security considerations.
- 1075 • The conference expects that many papers will be foundational research and not tied
1076 to particular applications, let alone deployments. However, if there is a direct path to
1077 any negative applications, the authors should point it out. For example, it is legitimate
1078 to point out that an improvement in the quality of generative models could be used to
1079 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1080 that a generic algorithm for optimizing neural networks could enable people to train
1081 models that generate Deepfakes faster.
- 1082 • The authors should consider possible harms that could arise when the technology is
1083 being used as intended and functioning correctly, harms that could arise when the
1084 technology is being used as intended but gives incorrect results, and harms following
1085 from (intentional or unintentional) misuse of the technology.
- 1086 • If there are negative societal impacts, the authors could also discuss possible mitigation
1087 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1088 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1089 feedback over time, improving the efficiency and accessibility of ML).

1090 11. Safeguards

1091 Question: Does the paper describe safeguards that have been put in place for responsible
1092 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1093 image generators, or scraped datasets)?

1094 Answer: [NA]

1095 Justification: Our paper doesn't release models that have the potential to cause harm like
1096 image generators or language models.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: Yes, we have well cited the packages.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: The paper provides code.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

1148 Answer: [NA]
 1149 Justification: The paper does not involve crowdsourcing nor research with human subjects.
 1150 Guidelines:
 1151 • The answer NA means that the paper does not involve crowdsourcing nor research with
 1152 human subjects.
 1153 • Including this information in the supplemental material is fine, but if the main contribu-
 1154 tion of the paper involves human subjects, then as much detail as possible should be
 1155 included in the main paper.
 1156 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
 1157 or other labor should be paid at least the minimum wage in the country of the data
 1158 collector.

1159 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
 1160 **Subjects**
 1161 Question: Does the paper describe potential risks incurred by study participants, whether
 1162 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
 1163 approvals (or an equivalent approval/review based on the requirements of your country or
 1164 institution) were obtained?
 1165 Answer: [NA]
 1166 Justification: The paper does not involve crowdsourcing nor research with human subjects.
 1167 Guidelines:
 1168 • The answer NA means that the paper does not involve crowdsourcing nor research with
 1169 human subjects.
 1170 • Depending on the country in which research is conducted, IRB approval (or equivalent)
 1171 may be required for any human subjects research. If you obtained IRB approval, you
 1172 should clearly state this in the paper.
 1173 • We recognize that the procedures for this may vary significantly between institutions
 1174 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
 1175 guidelines for their institution.
 1176 • For initial submissions, do not include any information that would break anonymity (if
 1177 applicable), such as the institution conducting the review.