

# Equi-explanation Maps: Concise and Informative Global Summary Explanations

TANYA CHOWDHURY, University of Massachusetts Amherst, USA

RAZIEH RAHIMI, University of Massachusetts Amherst, USA

JAMES ALLAN, University of Massachusetts Amherst, USA

In this work, we propose to summarize the model logic of a blackbox in order to generate concise and informative global explanations. We propose equi-explanation maps, a new explanation data-structure that presents the region of interest as a union of equi-explanation subspaces along with their explanation vectors. We then propose E-Map, a method to generate equi-explanation maps. We demonstrate the broad utility of our approach by generating equi-explanation maps for various binary classification models (Logistic Regression, SVM, MLP, and XGBoost) on the UCI Heart disease dataset and the Pima Indians diabetes dataset. Each subspace in our generated map is the union of  $d$ -dimensional hyper-cuboids which can be compactly represented for the sake of interpretability. For each of these subspaces we present linear explanations assigning weights to each explanation feature. We justify the use of equi-explanation maps in comparison to other global explanation methods by evaluating in terms of *interpretability*, *fidelity*, and *informativeness*. A user study further corroborates the use of equi-explanation maps to generate compact and informative global explanations.

Additional Key Words and Phrases: explainability, subspace interpretability, global explanations, explaining classifiers, model-logic subspaces

## ACM Reference Format:

Tanya Chowdhury, Razieh Rahimi, and James Allan. 2018. Equi-explanation Maps: Concise and Informative Global Summary Explanations. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

Wikipedia defines Explainable AI as *AI in which the results of the solution can be understood by humans*. Most models today accept a set of features (tabular or categorical) and combine them in a carefully constructed though often obscure way to produce a result. An “explanation” uses the same or different features to generate simple, interpretable information that gives an insight on how the model might have arrived at that result. For example, a complex neural model might be explained by a linear combination of a subset of the features. As machine learning models are increasingly being used in real-world decision making, it is important to provide explanations of model predictions to guide their use and to improve understanding of them.

Explanation algorithms which explain a single model prediction are known as *local* explanation algorithms, while those which approximate characteristics of an entire model are known as *global* explanation algorithms. Explanation algorithms which explain predictions by taking into consideration the original model parameters are known as *model-introspective* explainers. Methods which treat the original models as a black box, only to learn model characteristics

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 Association for Computing Machinery.

Manuscript submitted to ACM

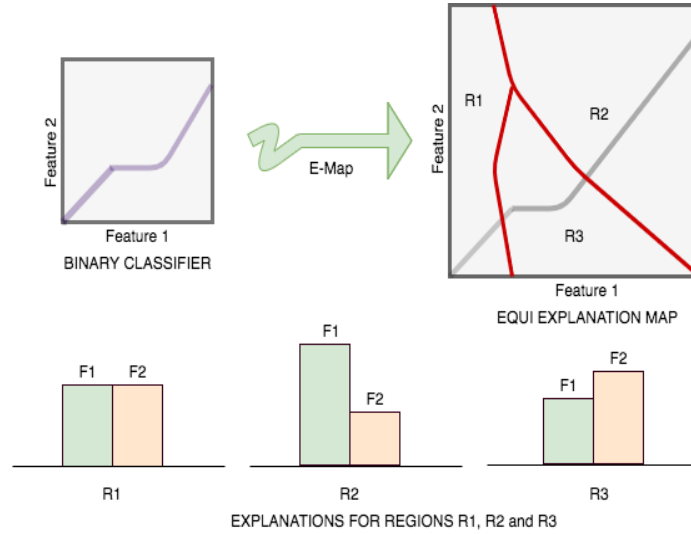


Fig. 1. A hypothetical binary classifier (top left) partitioned into equi-explanation maps using the proposed E-Map algorithm. The algorithm partitions the feature space (here containing 2 features  $f_1$  and  $f_2$ ) into three regions (R1, R2 and R3) based on similar model logic. Bar charts in the bottom row represent relative importance of the features in each of the 3 regions.

using secondary training data, are known as *model-agnostic* explainers. Algorithms in which we generate explanations for a model after it has been trained are known as *post-hoc* explanation algorithms. Explanation algorithms differ in *basic units* of explanation: some methods use the features as-is for interpretability, while some use mappings of features for the same. Different interpretability methods map the *interaction between features* to various degrees. Most explainability methods produce linear interpretations of model predictions, thus ignoring all inter-feature interaction terms. Generating explanations is also *time* and *resource* dependent. Some algorithms assign a *time budget*, and return the optimal explanation model derived within the assigned budget.

In this work, we focus on model-agnostic post-hoc linear interpretability techniques. The problem of *how the model logic varies across the input space* has not been studied well. Assume that a medical practitioner has to rely on an ML model decision to choose between different treatment plans for heart patients. Before relying on the model for such a critical decision, they would like a system to summarize the *basis* on which the model makes decisions for different values of patient statistics (e.g., *smoking* and *exercise*) [6]. Existing global linear explanation methods at best return a set of representative instances which cannot be used to give answers to such questions.

To generate more informative global explanations, we propose dividing the desired region of explanation features into *subspaces* based on similar logic, i.e.  $\epsilon$ -*equi-explanation* subspaces (Figure 1). In this work, we focus on the task of binary classification. Each  $\epsilon$ -equi-explanation subspace is a union of non-overlapping hyper-cuboids, each hyper-cuboid representing the range of values it covers over the explanation features. We employ a divide and conquer based approach, where in each step of the memoized recursion function, we compute the explanation vectors for each vertex of the obtained hyper-cuboid. We compute explanation vectors for instances based on their nearest Decision Boundary Point (DBP), a method which approximates LIME [12] results but with less uncertainty. We finally generate linear explanations for each equi-explanation region by aggregating its member hyper-cuboid explanation vectors based on weight. *In order to study a given blackbox model, our algorithm delivers compact, informative, summary global explanations presenting a set of equi-explanation regions and their corresponding explanation vectors.* We adapt relevant global explanation methods

such as SP-LIME[12], Guided-LIME[15], SHAP[10, 11], and MUSE[6] to form strong baselines and justify the usage of equi-explanation maps on grounds of *Interpretability*, *Fidelity*, and *Informativeness*. We also conduct a user study to demonstrate the effectiveness of our new explanation format in comparison to non-linear global explanation methods and show that equi-explanation maps outperform the strongest baseline by 40% on Fidelity and 36% on Informativeness on an average. In the spirit of reproducibility, our implementation is attached as supplementary material.

Our main contributions are as follows:

- We propose equi-explanation maps: a concise, informative new data structure to summarize the model logic of a blackbox in order to generate concise and informative global explanations. In doing so we propose the task of *Global Summary Explanation Generation*.
- We discuss different methods to generate equi-explanation maps using existing explanation algorithms (Appendix) and propose the E-Map architecture specifically for this task. We compare different methods to generate equi-explanations maps and show that E-Map outperforms studied algorithms by strong margins (Appendix).
- We propose new metrics to uniformly compare E-Map generated equi-explanation maps with other linear, additive global explanation methods. We also conduct a user study to prove the effectiveness of our proposed data-structure against mimic model based global explanation methods.

Once this work is accepted, all code and data required to replicate the experiments will be released in the spirit of reproducibility.

## 2 BACKGROUND

LIME [12] is a popular local explanation model in the machine learning literature. It is a model-agnostic linear explanation method that locally approximates a classifier with an interpretable model. It does so by perturbing inputs in a locality near the instance of interest, and generating labels for the perturbed inputs using the original model. Let  $\mathcal{G}$  be the class of interpretable models. The explanation model obtained by LIME, for an instance  $x$  is:

$$E(x) = \arg \min_{g \in \mathcal{G}} \mathcal{L}(f, g, \pi_x) + \Omega(g) \quad (1)$$

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in Z} \pi_x(z) (f(z) - g(z'))^2 \quad (2)$$

$$\pi_x(z) = \exp(-D(x, z)^2 / \sigma^2), \quad (3)$$

where  $\Omega(g)$  represents the complexity of the learnt explanation model,  $\pi_x$  describes the locality around  $x$  (usually represented by an exponential kernel on a distance function), and  $\mathcal{L}(f, g, \pi_x)$  is a model of how unfaithful the explanation model  $g$  is in neighbourhood of  $x$ .  $D$  denotes a distance function (cosine, here). In order to learn the local behaviour of  $f$  around  $x$ , they sample points uniformly at random in the proximity of  $x$ . These samples are generated by perturbing the original input, and replacing some feature values by zero. Rebeiro et al. [12] also propose SP-LIME, a method that selects a set of representative instances for global interpretability via sub modular optimization.

Laugel et al. [8] improve on LIME’s sampling techniques to optimize model fidelity. They generate surrogate-based explanations for individual predictions based on sampling centered on a particular relevant place of the decision boundary, rather than on the prediction itself. This allows them to achieve substantially better results, demonstrated visually on the UCI half moons datasets, where the local explanation often does not agree with the global explanation due to an unusually shaped decision boundary. Garreau et al. [3] derive closed form-expressions for linear explanation systems and report that the coefficients are proportional to the gradient of the function being explained. Reiger et al.

[13] present evidence that aggregate explanations are more robust to attacks than individual explanation methods. We use existing work [3, 7, 11, 12] along with concepts from linear algebra to propose a concise, informative new format for global explanation representation which can be employed in places where a condensed version of how model logic varies across the region of interest is a requirement.

### 3 EQUI-EXPLANATION MAPS

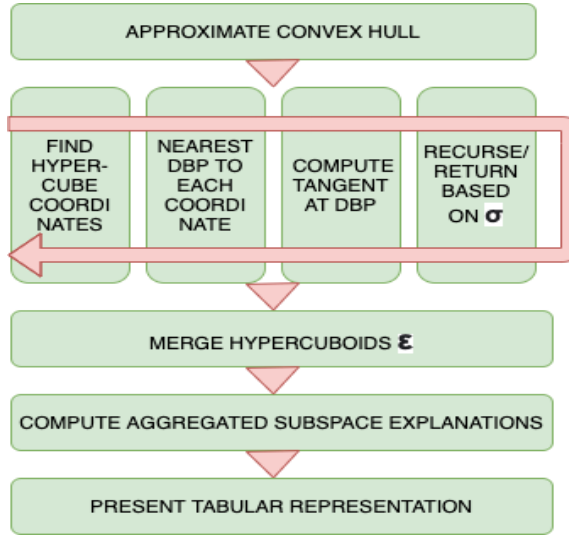


Fig. 2. E-map architecture to generate  $\epsilon$ -equi-explanation maps. Given a binary classifier  $f$  as a blackbox and region of interest  $P$ , we try to partition  $\mathcal{E}$  into equi-explanation regions in accordance to the decision boundary of  $f$ . We compute the approximate convex hull of  $\mathcal{E}$  and run a divide and conquer based approach, where in each recursion step we compute the explanation vectors of hyper-cuboid coordinates. Next, based on the standard deviation of the hyper-cuboid vertices exceeding  $\epsilon$  we decide if we want to divide the hyper-cuboid into sub hyper-cuboids. After the recursion ends, we try to merge hypercuboids, if their standard deviation is less than  $\epsilon$ . We lastly do a weighted aggregate of hyper-cuboid explanations to assign an explanation for a subspace.

Let  $f$  be a black-box binary classifier that maps input features  $\mathcal{S} = \{x_1, x_2, \dots, x_n\}$  to  $\{0, 1\}$ . The input features can be either tabular or categorical. Following other explanation models [12], the domain of an explanation model is the subset of interpretable features  $\mathcal{E}$  where  $\mathcal{E} \subseteq \mathcal{S}$  and  $|\mathcal{E}| = d$ . Assume that we have knowledge of a finite region of interest of  $\mathcal{E}$ , named  $P$ . The region of interest for explanation purposes may be smaller than the actual range of feature values. For example, healthcare providers would want to study heart disease symptoms in patients with age ranging from 0 to 120, instead of ages ranging from  $-\infty$  to  $+\infty$ .

We globally explain  $f$  by providing a division of region  $P$  into  $\epsilon$ -equi-explanation subspaces. An  $\epsilon$ -equi-explanation subspace is defined as a subspace of the explanation space where the deviations of local explanations for points within that subspace, do not exceed  $\epsilon$ . Specifically, classifier  $f$  is explained as  $[\mathbf{e}, \mathbf{W}]$  where  $\mathbf{e}$  is a partition of the explanation feature hyperspace  $\mathcal{E}$  into  $\epsilon$ -equi-explanation subspaces and each subspace  $e_i \in \mathbf{e}$  is explained with a function  $\mathbf{w}_i \in \mathbf{W}$ . Each explanation function  $\mathbf{w}$  belongs to a class of potentially interpretable models, for which we consider linear functions in this study. Thus, the explanation for each subspace  $\mathbf{w}_i \in \mathbb{R}^d$  represents the contribution of each explanation feature towards the model decision within that subspace. Standard deviation of linear explanation vectors is measured to check whether a subspace is  $\epsilon$ -equi-explanation or not.

To obtain  $\epsilon$ -equi-explanation subspaces, we propose a divide-and-conquer approach, **E-Map**, summarized by the pseudocode given in Algorithm 1. E-Map computes the hypercuboid of desired ranges of explanation features  $P$  and divides it into sub-hypercuboids if it is not an  $\epsilon$ -equi-explanation space. The obtained sub-hypercuboids are then recursively checked and divided if necessary. When all the obtained hypercuboids are  $\epsilon$ -equi-explanation subspaces, E-Map checks if neighboring subspaces can be merged to reduce the number of partitions. Each explanation subspace  $e$  thus consists of a set of hyper-cuboids in the  $d$ -dimensional space, and is linearly explained by a weighted average on each constituent hyper-cuboid's explanation vector. We describe the E-map approach in detail below.

---

**Algorithm 1** Pseudo-code of E-map approach to generate  $\epsilon$ -equi-explanation maps
 

---

**Require:**  $f$ : binary classifier,  $\mathcal{E} \in \mathbb{R}^d$ : explanation features,  $P$ : Region of interest for explanation,  $\epsilon \in [0, 1]$

```

1:  $C = \text{CONVEXHULL}(P)$ 
2:  $\text{CPoints} = \text{VERTICES}(C)$ 
3: procedure  $\text{DIVIDE-HYPER-CUBOID}(\text{CPoints})$ 
4:    $\text{DBP} = \text{DECISIONBOUNDARYPOINTS}(\text{CPoints})$ 
5:    $\text{ExplanationVectors} = \text{DBPTANGENTS}(\text{DBP})$ 
6:    $\sigma = \text{FINDDEVIATION}(\text{ExplanationVectors})$ 
7:   if  $\sigma > \epsilon$  then
8:      $\text{list} = \text{CREATESUBCUBOIDS}(\text{CPoints})$ 
9:     for  $\text{CPoint}$  in  $\text{list}$  do
10:       $\text{DIVIDE-HYPER-CUBOID}(\text{CPoint})$ 
11:   return  $\text{CPoints}, \text{ExplanationVectors}$ 
12:  $\text{CPoints}, \text{ExplanationVectors} = \text{DIVIDE-HYPER-CUBOID}(\text{CPoints})$ 
13:  $\text{CPoints}, \text{ExplanationVectors} = \text{MERGE-CUBOIDS}(\text{CPoints}, \text{ExplanationVectors})$ 
14:  $e, W = \text{AGGREGATED-SUBSPACE-EXPLANATION}(\text{Cpoints}, \text{ExplanationVectors})$ 

```

---

**Convex hull of region of interest.** The first step of E-Map is to compute the convex hull of our  $d$ -dimensional region of interest  $P$ . Finding an exact convex hull is an NP-hard problem, however, an approximate convex hull of the region is sufficient for our purpose. As a result, we compute the hull coordinates, using the maximum and minimum values of each feature in the region of interest, and get a hyper-cuboid hull of  $2^d$  vertices.

**Divide and conquer algorithm.** Starting with the hyper-cuboid hull of the region of interest, we use a divide-and-conquer algorithm to obtain an  $\epsilon$ -equi-explanation map. Given a hyper-cuboid hull, the function `divide-hyper-cuboid` in Algorithm 1 computes a local explanation of classifier  $f$  at each vertex of the hull. It then determines if the explanation vectors are similar enough to be labelled  $\epsilon$ -equi-explanation. If the explanation vectors are similar enough, the hyper-cuboid is an equi-explanation subspace. Otherwise, the hyper-cuboid is partitioned into two sub-hyper-cuboids. Splitting along the hyper-cuboid plain with the most explanation distance between its faces in the middle. Then, the function `divide-hyper-cuboid` is called recursively for each of the partitions.

**Generating local explanations.** Approaches like LIME use a falling exponential kernel defined on a distance metric to weigh the importance of each perturbed input as in Equation (3). On examining Equations (2) and (3), we observe that the decision boundary point (DBP) nearest to the instance being explained plays the most significant role in generating explanations. Points further than the nearest DBP will have little influence due to the rapidly falling negative exponential function. However, as Laugel et al. [7] point out, the solution by LIME largely depends on the density of sampling and the proximity being sampled. Considering the scale of our task, a method that generates approximate explanations but with less variability, would yield better results. Therefore, we compute the tangent to the

**Algorithm 2** Growing spheres algorithm for computing Decision Boundary Point nearest to  $x$  [7]**Require:**  $f : \mathbf{x} \rightarrow \{0, 1\}$ : a binary classifier,  $x \in \mathbf{x}$ : an observation to be explained,  $\eta, n$ : Hyperparameters**Ensure:** Nearest Decision Boundary Point  $e$ 

- 1: Generate  $(z_i)_{i \leq n}$  uniformly in  $SL(x, 0, \eta)$
- 2: **while**  $\exists k \in (z_i)_{i \leq n} : f(k) \neq f(x)$  **do**
- 3:      $\eta = \eta/2$
- 4:     Generate  $(z_i)_{i \leq n}$  uniformly in  $SL(x, 0, \eta)$
- 5:  $a_0 = \eta, a_1 = 2\eta$
- 6: Generate  $(z_i)_{i \leq n}$  uniformly in  $SL(x, a_0, a_1)$
- 7: **while not**  $\exists k \in (z_i)_{i \leq n} : f(k) \neq f(x)$  **do**
- 8:      $a_0 = a_1$
- 9:      $a_1 = a_1 + \eta$
- 10:     Generate  $(z_i)_{i \leq n}$  uniformly in  $SL(x, a_0, a_1)$
- 11: **return**  $k$ , the  $l_2$ -closest decision boundary point from  $x$

DBP nearest to the instance to be explained for obtaining the local linear explanation. This method leads to a solution approximate to what is returned by LIME, if the LIME instance neighbourhood has been thoroughly sampled. However, with a higher number of dimensions dense sampling around an instance is hardly practical. The proposed approach is thus expected to return solutions better than LIME with poor sampling.

**Nearest DBP<sub>x</sub>.** To find the point nearest to an instance  $x$  on the decision boundary of classifier  $f$ , we tweak the Growing Spheres algorithm proposed by Laugel et al. [7]. We uniformly sample points in the sphere centered at  $x$  with the radius of  $\eta$ , which is initially set to a large value. In order to sample points inside  $SL(x, a_0, a_1)$ , we sample observations uniformly distributed over the surface of a unit sphere, then draw  $\mathcal{U}(a_0, a_1)$ -distributed values and use them to re-scale the distances between the sampled observations and  $x$ . In order to sample uniformly on a unit sphere, we sample observations from  $\mathcal{N}(x, 1)$  and scale them to a distance of 1 from  $x$ . Any of the sampled points having a class different than  $x$  guarantees the presence of at least one DBP within the sphere. However, this may not be the nearest DBP to  $x$ . As a result, the process is repeated by sampling points in a spheres of radius  $\eta = \eta/2$  and keep halving  $\eta$  until we have a sphere where all sampled points have the same class label as  $x$ . This indicates the absence of a decision boundary in the sphere. Next we sample points inside the spherical layer  $SL(x, \eta, 2\eta)$  in search of points which have a class label different from  $x$ . Out of the sampled points which have a different label, we accept the one with the minimum  $L_2$ -distance to  $x$  as the nearest DBP. If no such point is found within this spherical layer, then we search in the next spherical layer  $SL(x, 2\eta, 3\eta)$ , and so on until the nearest DBP is found. Let us name this point  $DBP_x$ . The pseudo code for this process by [7] is provided in Algorithm 2.

**Tangent at DBP<sub>x</sub>.** Finding  $DBP_x$ , we then compute the tangent of the decision boundary of  $f$  at this point. To compute this, we randomly perturb point  $DBP_x$  and feed the perturbed instances to  $f$  to get their predicted labels. A weighted linear regression model is then learned on perturbed instances to obtain the tangent of  $f$  at  $DBP_x$ . We use the kernel function from KernelSHAP[11] as the distance metrics in this regression (Equation 4) along with the Loss function defined in Equation 2. This results in a  $d$ -dimensional tangent that locally explains instance  $x$ .

$$\pi_x(z) = \frac{M - 1}{\binom{M}{C_{|z|}} |z| (M - |z|)} \quad (4)$$

where  $|z|$  is the number of non-zero elements in the  $M$  dimensional vector  $z$ .

313  $\epsilon$ -**equi explanation subspaces**. We define the standard deviation  $\sigma_{EV}$  for a set  $EV$  of explanation vectors (Equation  
314 5) as:

$$315 \quad \sigma_{EV} = \sum_{\forall p_i, p_j \in EV} \sum_{m \in M} (p_{im} - p_{jm})^2 \quad (5)$$

318 We next compute  $\sigma_{EV}$ , where  $EV$  stands for the set of explanation vectors of a hyper-cuboid's coordinates. If this  
319  $\sigma$  exceeds the value of the hyperparameter  $\epsilon$ , the obtained region cannot be termed an  $\epsilon$ -equi-explanation subspace  
320 and is further divided. For further division, we iterate over each explanation feature ( $i \in M$ ) and find the average  
321 explanation vector corresponding to its minimum and maximum values ( $X_{i_{min}}$  and  $X_{i_{max}}$ ). For every feature  $i \in M$ , we  
322 then compute the  $L2$  distance between  $X_{i_{min}}$  and  $X_{i_{max}}$ . We pick the feature ( $f$ ) with the highest distance (Equation 6),  
323 and partition the hyper-cuboid into two across that feature.  
324  
325

$$326 \quad f = \arg \max_i (X_{i_{min}} - X_{i_{max}})^2 \quad (6)$$

328 If  $\sigma$  does not exceed  $\epsilon$ , an  $\epsilon$ -equi-explanation subspace has been found. When the recursion ends, a set of hyper-  
329 cuboids are obtained where the standard deviation of each hyper-cuboid is less than  $\epsilon$ . At this step, a merge function  
330 is used to check if any two hyper-cuboids can be merged while still satisfying the  $\epsilon$  constraint. These merged-hyper  
331 cuboids represent subspaces  $e$  of equi-explanation maps.  
332

333 **Subspace linear explanation.** Once we have obtained  $\epsilon$ -equi-explanation subspaces, each being a set of hyper-  
334 cuboids, we generate a linear explanation for each subspace depicting the behavior of  $f$  in that subspace. For this step,  
335 we first compute an explanation vector for each hyper-cuboid by computing the explanation vector for each of the  
336 hyper-cuboid's coordinates and averaging them. We additionally compute the volume of each hyper-cuboid (product of  
337 edge lengths). To generate an aggregated explanation for each subspace, we average each hyper-cuboid's explanation  
338 vector weighted by its volume. This leads to a division of the input space into  $\epsilon$ -equi-explanation maps with normalized  
339 explanation vectors for each subspace.  
340  
341

342 **Presentation of results.** E-Map partitions the space into subspaces where each subspace is a union of  $d$ -dimensional  
343 hyper-cuboids with  $d$  being the number of explanation features. For models with fewer than 3 explanation features,  
344 equi-explanation maps can be visualized as in Figure 1. However for models with greater than 3 explanation features,  
345 we present the coordinates of each hyper-cuboid in a compact tabular representation using  $2d$  numbers for each  
346 hyper-cuboid, to enable at-a-glance summaries (as shown in Table 1). The value of  $\epsilon$  can be set according to the  
347 granularity of explanations required.  
348  
349

## 350 4 EXPERIMENTS

352 To the best of our knowledge, our work is a first towards generating explanations that summarize the model logic of  
353 a blackbox. However since it lies in the large space of global explanation methods, we compare the gains of E-map  
354 generated equi-explanation maps with other global explanation methods. We also study other methods to generate  
355 equi-explanation maps and evaluate their quality with respect to E-map (Appendix).  
356  
357

### 358 4.1 Baseline Models

360 **SP-LIME** [12], an extension of LIME, is a model-agnostic approach proposed to choose diverse and representative  
361 instances to describe the global model logic. Given a budget  $B$ , SP-LIME selects  $|B|$  instances from a uniformly sampled  
362 set  $X$  using a greedy approach based on the local explanation for each instance.  
363  
364

Table 1. A compact equi-explanation map representation of a classifier trained on the UCI Heart disease dataset, using three explanation features: Age, RestECG and Cholesterol. Here we depict the three similar explanation regions (subspaces) proposed by our algorithm and present the approximate model logic corresponding to each of the three regions.

Subspace		Age	Cholesterol	RestECG	Explanation
1	Min	29	285	0	[0.26, <b>0.68</b> ,0.06]
	Max	44	364	2	
2	Min	44	126	0	[0.21,0.25, <b>0.54</b> ]
	Max	63	285	0	
	Min	29	364	1	
	Max	44	564	2	
	Min	63	284	2	
	Max	77	364	2	
3	Min	63	126	0	[ <b>0.45</b> ,0.2,0.35]
	Max	77	284	1	
4	Min	29	126	0	[0.33,0.31, <b>0.36</b> ]
	Max	44	285	1	
	Min	44	285	1	
	Max	63	364	2	
	Min	63	364	2	
	Max	77	564	2	

Table 2. Accuracy of the four chosen classifiers on the training and test sets of the Heart disease and Pima Indians dataset.

Dataset	Algorithm	Training Accuracy	Test Accuracy
Heart Disease	Logistic Regression	0.86	0.80
	SVM	0.92	0.80
	MLP	1	0.81
	XGBoost Classifier	1	0.78
Pima Diabetes	Logistic Regression	0.79	0.73
	SVM	0.83	0.72
	MLP	1	0.77
	XGBoost Classifier	1	0.75

**Guided-LIME** [15] adds a structured-sampling preprocessing step to the input of SP-LIME in order to improve the fidelity of LIME-based approaches. To do this, they employ Formal Concept Analysis (FCA) assuming access to the complete model training data. Generating SP-LIME explanations on the full dataset (especially for tabular features) has large complexity, which Guided-LIME successfully reduces.

**SHAP** [11], originally a local feature attribution method, is extended for global explanation of tree-based models [10]. The authors report that it outperforms existing explainers on various metrics like run time, accuracy, consistency guarantees, mask, resample, and impute for tree based models.

**MUSE** [6] is a rule-based mimic model [1] based explanation algorithm to explain how a model behaves in subspaces characterized by certain features of interest. It aims to learn compact *decision sets*, each of which is a series of if-then rules built by optimizing for fidelity, unambiguity, and interpretability. Since MUSE and equi-explanation maps have representational differences, we compare equi-explanation maps to MUSE only by the user study. We do not include comparison with Interpretable decision sets (IDS [5]) and Bayesian decision lists (BDL) [9], as they have been shown to under-perform MUSE [6].



## 4.2 Evaluation Metrics

To enable comparison of equi-explanation maps to other linear global explanation algorithms, each with a different representation format, we propose the following general evaluation metrics.

**Interpretability:** The multiplicative inverse of the amount of information (numbers) needed to present generated explanations to users. It is dependent on the format of explanation presentation by an explanation algorithm. For example, if a model with 10 input features and 4 explanation features is to be explained with 3 representative instances, interpretability would be  $\frac{1}{3*(10+4)}$ . The higher the score, the more interpretable the explanation algorithm.

**Fidelity:** The fraction of sampled points from the region of interest for which the blackbox prediction agrees with the *explanation model prediction*. For explanation models that generate representative instances, the reconstructed *explanation model prediction* is either the prediction by the explanation vector of its subspace (if subspaces are defined) or the prediction by the explanation vector of its nearest representative element. The higher the fidelity, the better the explanation algorithm.

**Informativeness:** The average similarity between local and subspace explanation vectors for points uniformly sampled in the region of interest. The local explanation of a sampled instance is computed using LIME. The subspace explanation is either the explanation vector of the subspace (if subspaces are defined) or the explanation vector of the representative unit nearest to the instance. Similarity is computed using cosine similarity. Higher this metric, more informative the summary.

## 4.3 Experimental Settings

**Datasets :** We perform our experiments on two real-world datasets from the medical domain. The motivation is the known importance of subspace explanations in clinical diagnostic settings [5]. The first dataset is the *UCI Heart Disease dataset*<sup>1</sup>, which has 303 instances and 14 real valued features designed to predict the presence (labels: 1,2,3,4) and absence (label: 0) of heart disease. Secondly we use the *Pima Indians Diabetes dataset*<sup>2</sup>, predicting the onset of diabetes within 5 years in Pima Indians, given their medical details. It is a binary classification dataset with 768 observations containing 8 input features and a binary output. The label for each instance is either 0 or 1, with 1 indicating that the person would see an onset of diabetes within 5 years.

**Classifiers :** We train binary classifier models on the Heart Disease and Pima Diabetes datasets using four algorithms Logistic Regression, SVM, MLP, and XGBoost. We use the scikit-learn implementation of logistic regression. We learn a Support Vector Machine with an RBF kernel using the scikit-learn implementation. We experiment with different configurations of MLP using PyTorch. We settle on an architecture with 3 hidden layers. For a dataset of  $N$  dimensions, the MLP has  $N$  neurons in the first layers,  $2 * N$  in the second,  $N$  in the third,  $\frac{N}{2}$  in the fourth and a single neuron in the final output layer. For the XGBoost classifier, we use `XGBClassifier` from `xgboost`. The prediction threshold is set to 0.5, i.e., the prediction is 1 if scores are greater or equal to 0.5, otherwise the prediction is zero[20]. We split both datasets into three parts: train, validation, and test in the ratio of 80:10:10. We use five-fold cross validation on the training data to learn the supervised models. We compute the classification accuracy of each of these algorithms and report results in Table 2.

**Setting 1 :** In the first set of experiments, we compare the performance of Equi-explanation maps with respect to other global linear explanation methods on the two chosen datasets. We carry out experiments with all input features as explanation features and with a budget of four representative instances (as in SP-LIME). In order to carry out fair

<sup>1</sup><https://archive.ics.uci.edu/ml/datasets/Heart+Disease>

<sup>2</sup><https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.csv>

Table 3. Comparing E-map generated equi-explanation maps to existing global explanation algorithms with a budget of 4 representative instances on all explanation features for two medical domain datasets, averaging performance over all four classifiers.

Method	UCI Heart Disease			PIMA Indian Diabetes		
	Interpretability	Fidelity	Informativeness	Interpretability	Fidelity	Informativeness
SP-LIME	0.083	0.49	0.56	0.083	0.56	0.63
Guided-LIME	0.083	0.62	0.65	0.083	0.62	0.66
SHAP	0.083	0.56	0.66	0.083	0.60	0.69
<b>E-Map</b>	0.017	<b>0.86</b>	<b>0.79</b>	0.017	<b>0.88</b>	<b>0.86</b>

Table 4. Comparing E-map generated equi-explanation maps to explain four different classifiers on two medical domain datasets for a fixed value of  $\epsilon$  set to 0.6.

Method	UCI Heart Disease			PIMA Indian Diabetes		
	Interpretability	Fidelity	Informativeness	Interpretability	Fidelity	Informativeness
Logistic Regression	0.037	0.97	0.91	0.050	0.99	0.95
SVM	0.018	0.56	0.65	0.025	0.69	0.82
3-layer MLP	0.032	0.82	0.82	0.040	0.87	0.93
XGBoost Classifier	0.021	0.80	0.79	0.028	0.93	0.87

comparison, we tune E-map with different values of  $\epsilon$  until four subspaces are generated. For the sake of uniformity, we also demarcate the centroid of each subspace in E-map as its representative vector. SP-LIME and Guided-LIME algorithms output representative instances and their respective explanations as output, but no subspace information. Equi-explanation maps output representative instances, their explanations, and the coordinates of subspace hyper-cuboids (refer to Table 1). Both SP-LIME and Guided-LIME require a sampling density to determine the number of perturbations, the value of which is retained from the original LIME repository. For E-map, we set initial sampling radius  $\eta$  to 1 and the number of points to sample on the sphere  $n$  to 1,000 (following recommendations by Laugel et al. [7]). For the comparing methods, we set as many parameters as possible to the values reported in their original drafts or repositories. MUSE is a mimic model-based explanation method whose performance is difficult to compare with representation vector-based explanation empirically. As a result, we put off comparison of E-map with MUSE to a user study (Section 5). For our evaluation metrics *Fidelity* and *Informativeness*, we uniformly sample 500 points from the region of interest. The same 500 points are used to evaluate all competing systems for a given data set and classifier.

**Setting 2 :** For the next set of experiments, we compare the equi-explanation maps generated by E-Map for the four different classifiers described above. Setting the value of  $\epsilon$  to a fixed 0.6, the aim of this experiment is to study the variation in *Interpretability*, *Fidelity* and *Informativeness* for explanations of classifiers with different complexities on the two chosen datasets. As above, we set initial sampling radius  $\eta$  to 1 and the number of points to sample on the sphere  $n$  to 1,000 (as above). We again sample 500 points from the region of interest and use the same points to evaluate explanations for different classifiers.

#### 4.4 Results and Observations

The results of our experiments comparing equi-explanation maps generated by E-Map with representation-vector based explanations by other global explanation methods are reported in Table 3. The representative instances returned by equi-explanation maps show 43% and 38% higher fidelity than those returned by SP-LIME and Guided-LIME, respectively. This indicates that a user is more likely to guess the blackbox model’s decision for an instance, when shown an E-map

521 explanation as compared to when shown a \*-LIME explanation. This might be attributed to equi-explanation maps  
522 presenting subspaces of complex shapes and sizes considering intricacies of decision boundaries compared to the  
523 spherical subspaces carved by \*-LIME. E-Map explanations are also 41% and 21.5% more informative as compared to  
524 SP-LIME and Guided-LIME explanations. This indicates that a user is more likely to accurately guess the explanation of  
525 an unseen instance when shown an equi-explanation map explanation rather than the when shown other kinds of  
526 explanations. This might again be attributed to complex subspace boundaries for equi-explanation regions as compared  
527 to other explanation algorithms.  
528

529 Due to the extra reporting of subspace coordinates, equi-explanation maps show 79.6% lesser interpretability  
530 compared to the other approaches. Since SP-LIME and Guided-LIME only return representative instances and their  
531 explanations, their interpretability score for a budget B is computed as  $|B| * (\text{NUMBER OF INPUT FEATURES} + \text{NUMBER OF}$   
532  $\text{EXPLANATION FEATURES})$ . The interpretability of E-map, which returns hyper-cuboid dimensions, additionally includes  
533 an extra  $|B| * (2 * \text{NUMBER OF EXPLANATION FEATURES} * \text{NUMBER OF HYPER-CUBOIDS})$  term. As a result, although we see a  
534 gain in Fidelity and Informativeness with E-map, we see a drop in Interpretability.  
535

536 The results of comparing E-map explanations for different classifiers is presented in Table 4. Over experiments  
537 with E-map on different classifiers, we observe a strong correlation between the number of subspaces returned by  
538 E-Map and the complexity of the classifier being explained, for fixed values of  $\epsilon$ . Overall the simplest classifier (logistic  
539 regression) achieves 105% more interpretability, 79% more fidelity and 40% more on informativeness as compared to  
540 equi-explanation maps' explanation of the classifier with the most complex decision boundary (here: SVM) for  $\epsilon = 0.6$ .  
541 This is intuitive as the more curved the decision boundary is, the greater the deviation in explanations is and vice versa.  
542

543 MUSE is a series of cascading if-then rules, with a maximum hierarchy of two layers. For datasets with mostly  
544 tabular features, MUSE will have exponential combinations over feature ranges and would be pretty unintuitive to  
545 observe. Decision set-based approaches seem to work well for datasets with mostly categorical features, like the datasets  
546 chosen in demonstrated examples of previous studies [5, 6, 9]. Since many real world datasets are largely tabular,  
547 we believe equi-explanation maps would provide better explanations for them as compared to decision sets-based  
548 approaches [5, 6, 9].  
549  
550  
551

## 552 5 USER STUDY

553 In the previous section, we saw that our approach outperforms the compared baselines on grounds of Fidelity and  
554 Informativeness. However different global explanation algorithms result in different formats, making it hard to compare  
555 algorithm performances. In order to enable fair comparison between rule-list based mimic model explanations (MUSE)  
556 and equi-explanation maps, we conduct an user study. Inspired by explainable AI literature, we recruited 20 students  
557 who had completed at least one graduate-level machine learning course. These students were divided into two equal  
558 sized groups and each group was presented with two types of multi-choice questions (MCQ). Group A was shown  
559 equi-explanation maps for an XGBoost classifier trained on the UCI Heart disease dataset while Group B was shown  
560 global decision sets (MUSE explanations) generated to explain the same classifier. The questions are of two types : (i)  
561 Given the global explanation, predict the output by the model for a new instance. (ii) Given an input instance and its  
562 prediction, indicate which two features are the most important in making that prediction by the classifier. The user's  
563 response accuracy and the time taken for each response were measured. For question 2, the users were additionally  
564 asked if they felt confident enough to make a prediction (Yes/No). The ground truth for question 1 was obtained by  
565 feeding the asked instance to the black box classifier. The ground truth for question 2 was obtained by generating LIME  
566 explanations for the asked instance and picking the two features with the highest weights.  
567  
568  
569  
570  
571  
572

We observed that MUSE achieves a mildly higher accuracy/time ratio for question 1, outperforming equi-explanation maps by 11%. While the accuracy of both explanation algorithms are comparable, the time taken by the MUSE group to predict the black box output is 26% less than the time required to predict the same for the equi-explanation maps group. However, the equi-explanation map group outperforms the MUSE group by 52% on accuracy/time ratio in question 2. The equi-explanation map group sees both higher accuracies (46%) and lower times (22%) than the MUSE group. Additionally, 70% of the equi-explanation map group in question 2 felt confident answering the questions as compared to only 30% of the users in the MUSE group.

This shows that if an user wants to understand what features influence the decision making in a certain region (which is the primary intent behind generating explanations), equi-explanation maps should be unarguably preferred.

## 6 DISCUSSION

Local explanations are useful to understand the blackbox intent for a specific instance, but do not tell much about the larger picture. Most existing global explanation methods on the other hand are too sparse and not very informative about the variation in model logic across the region of interest[7]. Global explanations that summarize the model logic using equi-explanation maps lie somewhere in between, with a informativeness-interpretability tradeoff, which can be set by tuning  $\epsilon$  as per the user requirements.

Most existing explanation work focuses solely on *Interpretability* and *Fidelity/Faithfulness* as the primary metrics to optimize during explanation generation. However we believe that *Informativeness* - which is a proxy for the knowledge gained is an important metrics to consider as well. While the metric *Fidelity* is an evaluation metric in the blackbox decision space, *Informativeness* is an evaluation metric in the explanation space. We believe that lower bounds and confidence intervals on Fidelity and Informativeness should be mandated in instances where the explanation is highly consequential: e.g. deciding between treatment options by a doctor, deciding jail term length, etc. We believe that using summaries of the black box logic for global explanations instead of existing approaches will induce more trust in users in the above mentioned scenarios.

Due to the increased number of bits in Equi-explanation maps based explanations, they are more suitable in low dimensional settings. However they can still be generated in higher dimensional settings to verify black box model behavior. For example: a developer would prefer an informative summary explanation to verify how their model prioritizes features on all subspaces, even at the cost of interpretability.

Studying explanations with a subset of input features as explanation features might not always be as insightful even if it is more interpretable. For example, for a classifier with input features [A,B,C] if we generate explanations using only feature A and B, there might be a causal feature C driving model decision making for that variable. However, these factors depend on the exact problem statement in hand.

## 7 RELATED WORK

Recently Setzu et al. [18] propose GLocalX a tool that adds an interpretable layer on top of a blackbox by aggregating local explanations agnostic to the model being explained. GLocalX hierarchically aggregates the local explanations, represented as decision rules with the goal of emulating the blackbox. Their output format is similar to MuSe [6] and serves as a mimic model to the blackbox. Our proposed data structure, subspaces as a union of hyper-cuboids can also be visualized using Polyhedral Sets from linear algebra. Ruggieri et al. [14] propose a method of learning a parameterized linear system whose class of polyhedra includes a given set of example polyhedral sets and it is minimal.

625 Apart from explanations for black box models, certain algorithms generate explanations while taking the model  
626 architecture into consideration. One such notable model-introspective explanation method for deep learning models  
627 is DeepLIFT [19] which backpropagates the output of DNNs to assign each input feature a contribution weight.  
628 DeepLIFT [19] specializes in that its assigned feature weight, positive or negative, can be computed in a single backward  
629 pass of the neural network. Selvaraju et al. [17] introduce Grad-CAM, a method for generating *visual explanations*  
630 across the layers of a convolutional neural network (CNN), using target gradients to create coarse localization maps  
631 highlighting region importance. SHAP [11] provides unifying framework and formalizes explainability using the  
632 Shapely values principle from game theory. The SHAP method assigns each feature an importance value for a particular  
633 prediction. It is notable because it proves that there exists a unique solution in this class with a set of desirable properties.  
634 Their framework unifies six existing methods including LIME [12] and DeepLift [19] discussed above. Apart from the  
635 local explanation techniques discussed so far, there are a few algorithms that help users make global conclusions about  
636 the model.  
637

638  
639  
640 The above methods focus on features which are present, even though these features might have a negative contribution  
641 in the classification. A few recent works have been focused on identifying features which are necessary or sufficient  
642 to explain an instances classification by a model. Dhurandhar et al. [2] present contrastive explanations to explain  
643 black box classification models. Given an input, they find what features must be *minimally and sufficiently present* and  
644 *minimally and sufficiently absent* to justify its classification. The authors argue that this format of explanations is more  
645 in line with the human way of thinking. Another recent branch of model agnostic black box explanation methods  
646 include counterfactual explanations. Given a *query* image  $\mathcal{I}$ , for which a vision system predicts class  $C$ , a counterfactual  
647 visual explanation identifies how  $\mathcal{I}$  could change such that the system would output a different specified class [4].  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660

## 660 8 CONCLUSION

661 In this work, we proposed the new paradigm of summarizing the model logic of a blackbox in order to generate  
662 global explanations. In order to do this, we propose Equi-explanation maps, a novel concise representation for global  
663 explanations. We further proposed E-Map, an effective method that generates equi-explanation maps. Using hyper-  
664 cuboids as units of equi-explainability, we termed the union of hyper-cuboids to be a subspace and assigned a linear  
665 explanation to each subspace. We experimented on two medical records datasets: the UCI heart disease dataset and  
666 the Pima Indians Diabetes dataset and our approach was evaluated using the metrics interpretability, fidelity and  
667 informativeness, and substantially outperformed competitive methods in most of these. With this work, we introduce the  
668 task of *Global summary explanation generation* - explanations which present a summary of the model logic of a blackbox  
669 model. We hope future explainability researchers study the compactness-informativeness tradeoff for global summaries  
670 and propose better ways to generate summary explanations. It would also be more effective to see equi-explanation  
671 maps being generated using causal explanation techniques with bounds on subspace uncertainty.  
672  
673  
674  
675  
676

## 9 APPENDIX - COMPARING DIFFERENT METHODS TO GENERATE EQUI-EXPLANATION MAPS

To evaluate the performance of the proposed algorithm in generating  $\epsilon$ -equi-explanation maps, we adapt relevant algorithms in existing literature.

### 9.1 Baselines

The baseline methods we compare with for the problem statement are as follows:

**Random-Division.** This is a rudimentary baseline. We uniformly sample a small number of points in the space of the explanation features. We divide the region of interest into hyper-cuboids based on these points. Next, we compute the LIME explanation vectors of each of the coordinates of these hyper-cuboids. If the standard deviation of any of these hyper cubes exceeds  $\epsilon$ , we randomly sample points within the hypercube and divide it into sub-hypercubes. This recursive process of randomly driven division continues until the  $\epsilon$  requirement is met. Once the standard deviation of a hyper-cuboid is less than  $\epsilon$ , we run a merge operation to see if a particular cube can be merged to its neighbors in order to reduce the number of subspaces.

**Subspace-LIME.** LIME derives explanations by perturbing inputs in the vicinity of the instance to be explained and contributions are weighted by an exponentially falling kernel which is a function of distance. In this baseline we test the scalability of LIME for generating subspaces on the basis of similar model logic. We uniformly sample points in the convex hull of the region of interest ( $C$ ) and compute LIME explanations for each point. We then start merging these explanation vectors bottom up, as long as the standard deviation of the explanation vectors of points in a group does not exceed  $\epsilon$ . To facilitate this merging, we use hierarchical agglomerative clustering and use the Scipy implementation of the same. We then compute the min-max convex hull of each cluster to transform them into hypercuboids. The average explanation of all points in a cluster is assigned as the subspace's linear explanation.

### 9.2 Evaluation Metrics

We propose evaluation metrics to compare the above proposed systems for generating equi-explanation maps. Each of the competing systems explained above returns a set of subspaces each of which is a union of hypercubes. These subspaces are exhaustive, i.e., the union of these subspaces fully covers the input space  $C$ .

**Number of Subspaces.** The number of partitions in the  $\epsilon$ -equi-explanation map generated by a particular algorithm. The higher the number of subspaces for a fixed value of  $\epsilon$ , the less human readable the solution is, and hence lower the interpretability. Interpretability of an explanation, for this set of experiments is thus inversely proportional to Num-Subspaces for a given value of  $\epsilon$ .

**Faithfulness.** By clustering a set of points in the same subspace, an algorithm indicates that the model takes decisions based on similar logic on each of these points. Based on that, it proposes a representative explanation vector. We random-uniformly sample points from each subspace and use the subspace representative explanation vector to regenerate the blackbox labels. In this metric, we report the precision of the representative explanation vectors in regenerating model predictions within their subspaces.

### 9.3 Hyperparameter Settings

There are two main hyperparameters in the proposed approach:  $\epsilon$  and  $\eta_g$ . Hyperparameter  $\epsilon$  controls the variance between explanation vectors in a subspace. High values of  $\epsilon$  indicating high variance subspaces. There is also a trade-off between  $\epsilon$  and the number of subspaces generated by the algorithm.  $\eta_g$  impacts the granularity of the generated

Table 5. Comparing performance of competing models on proposed evaluation metrics: Number of Subspaces (Interpretability), Faithfulness, and Generation time for both datasets on different classifiers, with  $\epsilon = 0.5$ .

Method	Classifier	UCI Heart Disease		PIMA Indian Diabetes	
		N.Subspaces	Faithfulness	N.Subspaces	Faithfulness
Random-Division	Logistic Regression	16	0.78	14	0.82
Subspace-LIME		12	0.82	12	0.89
E-map		7	0.81	7	0.92
Random-Division	SVM	19	0.76	16	0.83
Subspace-LIME		14	0.81	14	0.89
E-map		10	0.83	11	0.91
Random-Division	MLP	21	0.67	19	0.82
Subspace-LIME		15	0.72	15	0.89
E-map		9	0.71	10	0.85
Random-Division	XGBoost	19	0.59	21	0.76
Subspace-LIME		12	0.64	18	0.82
E-map		11	0.66	12	0.85

subspaces. A very small  $\eta_g$  nearing zero, would learn perfectly curved surface boundaries and a high number of hyper-cuboid per subspace. Hence there is a trade-off between  $\eta_g$  and the human-readability of an explanation. Based on preliminary analysis, we set  $\eta_g$  to 1 based on the level of visualization detail offered by some platforms. We apply grid search over a range of values of  $\epsilon$  (0.1,1,5,10,20,30,50,100) for deciding on a value for  $\epsilon$  in our experiments. For Random-Division, we start the partitioning algorithm with 1 point. For generating explanations in Subspace-LIME, we set all parameters to default in the LIME API. For hierarchical clustering, we use the *linkage* method of scipy, where we set our metric to a custom function which returns distances in the explanation space.

## 9.4 Results

In our first set of experiments, we set the value of  $\epsilon$  to 0.5 and report the results in Table 5. We observe that our proposed approach, E-map, generates the least number of subspaces (30% fewer than Subspace LIME on the UCI Heart Disease dataset and 32% less than Subspace LIME on the PIMA Indians Diabetes dataset) making it the most interpretable approach. We see that as the decision boundary gets more complex, the number of subspaces and the generation time of all three algorithms increases. The faithfulness of explanations resulting from all three approaches are similar. However the union of all subspaces generated by Subspace-LIME does not cover  $C$ . Hence it is beneficial to use E-map for reliable decision making. We observe that E-map (which uses DBP sampling) produces explanations strongly correlated to explanations by LIME (correlation coefficient: 0.76).

## 10 ACKNOWLEDGEMENTS

This work was supported in part by the Center for Intelligent Information Retrieval, in part by NSF grant IIS-2039449 and in part by NSF grant number 1813662. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect those of the sponsor.

## REFERENCES

- [1] Nadia Burkart and Marco F Huber. A survey on the explainability of supervised machine learning. *Journal of Artificial Intelligence Research*, 70:245–317, 2021.
- [2] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. *arXiv preprint arXiv:1802.07623*, 2018.

- 781 [3] Damien Garreau and Ulrike Luxburg. Explaining the explainer: A first theoretical analysis of lime. In *International Conference on Artificial Intelligence*  
782 *and Statistics*, pages 1287–1296. PMLR, 2020.
- 783 [4] Yash Goyal, Ziyang Wu, Jan Ernst, Dhruv Batra, Devi Parikh, and Stefan Lee. Counterfactual visual explanations. In *ICML*, pages 2376–2384. PMLR,  
784 2019.
- 785 [5] Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. Interpretable decision sets: A joint framework for description and prediction. In  
786 *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1675–1684, 2016.
- 787 [6] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. Faithful and customizable explanations of black box models. In *Proceedings of*  
788 *the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 131–138, 2019.
- 789 [7] Thibault Laugel, Marie-Jeanne Lesot, Christophe Marsala, Xavier Renard, and Marcin Detyniecki. Comparison-based inverse classification for  
790 interpretability in machine learning. In *International Conference on Information Processing and Management of Uncertainty in Knowledge-Based*  
791 *Systems*, pages 100–111. Springer, 2018.
- 792 [8] Thibault Laugel, Xavier Renard, Marie-Jeanne Lesot, Christophe Marsala, and Marcin Detyniecki. Defining locality for surrogates in post-hoc  
793 interpretability. *ICML '18 Workshop on Health*, 2018.
- 794 [9] Benjamin Letham, Cynthia Rudin, Tyler H McCormick, David Madigan, et al. Interpretable classifiers using rules and bayesian analysis: Building a  
795 better stroke prediction model. *Annals of Applied Statistics*, 9(3):1350–1371, 2015.
- 796 [10] Scott M Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M Prutkin, Bala Nair, Ronit Katz, Jonathan Himmelfarb, Nisha Bansal, and  
797 Su-In Lee. From local explanations to global understanding with explainable ai for trees. *Nature machine intelligence*, 2(1):2522–5839, 2020.
- 798 [11] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in neural information processing systems*, pages  
799 4765–4774, 2017.
- 800 [12] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the*  
801 *22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- 802 [13] Laura Rieger and Lars Kai Hansen. Aggregating explanation methods for stable and robust explainability. *arXiv preprint arXiv:1903.00519*, 2019.
- 803 [14] Salvatore Ruggieri. Learning from polyhedral sets. In *Twenty-Third International Joint Conference on Artificial Intelligence*, 2013.
- 804 [15] Amit Sangroya, Mouli Rastogi, C Anantaram, and Lovekesh Vig. Guided-lime: Structured sampling based hybrid approach towards explaining  
805 blackbox machine learning models.
- 806 [16] Patrick Schwab and Walter Karlen. Cxplain: Causal explanations for model interpretation under uncertainty. *arXiv preprint arXiv:1910.12336*, 2019.
- 807 [17] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations  
808 from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- 809 [18] Mattia Setzu, Riccardo Guidotti, Anna Monreale, Franco Turini, Dino Pedreschi, and Fosca Giannotti. Glocalx-from local to global explanations of  
810 black box ai models. *Artificial Intelligence*, 294:103457, 2021.
- 811 [19] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings*  
812 *of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, page 3145–3153. JMLR.org, 2017.
- 813 [20] ShubhankarRawat. ShubhankarRawat/heart-disease-prediction. 2016. [Online; accessed 11-June-2021].
- 814 [21] Arnaud Van Looveren and Janis Klaise. Interpretable counterfactual explanations guided by prototypes. *arXiv preprint arXiv:1907.02584*, 2019.
- 815 [22] David Watson. Interpretable machine learning for genomics. 2021.