

GALE: GLOBALLY ASSESSING LOCAL EXPLANATIONS

Peter Xenopoulos & Claudio Silva

New York University
{xenopoulos, csilva}@nyu.edu

Gromit Chan

Adobe Research
ychan@adobe.com

Harish Doraiswamy

Microsoft Research India
harish.doraiswamy@microsoft.com

Luis Gustavo Nonato

Universidade de São Paulo
gnonato@icmc.usp.br

Brian Barr

Capital One
brian.barr@capitalone.com

ABSTRACT

Local explainability methods – those which seek to generate an explanation for each prediction – are increasingly prevalent. However, results from different local explainability methods are difficult to compare since they may be parameter-dependant, unstable due to sampling variability, or in various scales and dimensions. We propose GALE, a topology-based framework to extract a simplified representation from a set of local explanations. GALE models the relationship between the explanation space and model predictions to generate a topological skeleton, which we use to compare local explanation outputs. We demonstrate that GALE can not only reliably identify differences between explainability techniques but also provides stable representations. Then, we show how our framework can be used to identify appropriate parameters for local explainability methods. Our framework is simple, does not require complex optimizations, and can be broadly applied to most local explanation methods.

1 INTRODUCTION

Increasingly complex machine learning (ML) models are being deployed in industries such as healthcare, cybersecurity, and banking. While industries welcome the performance boost from these ML models, organizations are also starting to require their models to provide clear explanations for their predictions, as necessitated by regulations, such as GDPR. Thus, explainable artificial intelligence (XAI) techniques are especially relevant, particularly those which provide explanations for individual samples.

There are many techniques which provide *local* explanations for predictive models (Ribeiro et al., 2016). These techniques typically generate a real-valued vector which represents the attributions of the input features to the prediction. However, even using a single explanation technique, it is common to observe different outcomes based on different parameters, which can be hard to tune. Furthermore, since some local explainability methods rely on random sampling, there may also exist sampling variability, which can further impact the stability of the explanation output. Since there are many local

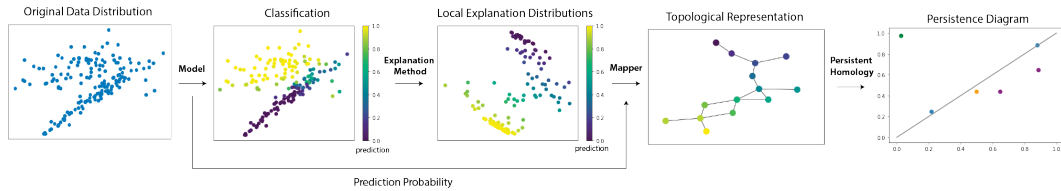


Figure 1: GALE transforms local explanations into a topological representation used to compare local explanation methods.

explanation methods, each with their own set of parameters, it raises an important question – *how do we assess and compare explanations from different local explainability frameworks?* By doing so, for a given prediction problem, we can develop a quantitative sense of “consensus” among various local explanation frameworks for different prediction tasks, which can be useful for benchmarking.

We propose GALE (globally assessing local explanations), a simple approach to assess the difference between sets of local explanations. We do so through an analysis of the topological properties of a given explanation space for binary classification problems. Specifically, we model an explanation method as a scalar function that captures the relationship between the explanation space and the predictions. Using this function, we compute its topological skeleton. This skeleton is used to generate a topological signature which is then used to compare explanation methods. We can calculate distances between signatures produced by different local explanation methods to compare their similarity. GALE is easy to implement and lacks complex optimizations or parameter tuning.

To the best of our knowledge, our approach is the first to use computational topology for comparing XAI methods. We view topology as a promising direction for understanding and comparing explanation methods and GALE, a powerful yet accessible framework, is the first step in that direction. The contributions of this work can be summarized as follows:

1. GALE (Global Assessing Local Explanations) a topology-based approach to generate a global signature for a given local explanation method’s output. By providing a domain-agnostic signature for explanation techniques, our approach allows comparison across heterogeneous explanation approaches.
2. We demonstrate that GALE is both stable and can elucidate differences between explanations through experiments on gradient-based and surrogate model techniques with real-world and synthetic data.
3. We show that GALE can be used to find appropriate parameters for local explanation methods by comparing topological signatures for different sets of parameters.

2 RELATED WORK

To balance the desire for machine learning model performance improvements along with growing calls for explainable model decisions, practitioners are increasingly turning to posthoc local explainability methods (Weld & Bansal, 2019). These local explanation methods typically produce a vector representing the attributions of input features, and are generally well-received by users. For example, LIME (Ribeiro et al., 2016) and SHAP (Lundberg & Lee, 2017) are two common posthoc, model-agnostic methods that produce local feature attributions. However, it is hard to compare results across different explainability methods. Jeyakumar et al. (2020) find, through surveys with hundreds of non-technical users, that explanation-by-example and LIME were the preferred explanation styles. Barr et al. (2020) propose synthetic data generation with known ground truth to explore explainability. Taxonomies, such as *Explainability Fact Sheets*, also provide some guidance on a unified language around explainability and raise important questions for explainability researchers to address (Doshi-Velez & Kim, 2017; Sokol & Flach, 2020).

Some local explanation methods, like integrated gradients, rely on baselines to generate feature attributions. Haug et al. (2021) show that popular local explanation methods are sensitive to the choice of baseline. Sundararajan & Najmi (2020) propose Baseline Shapley and prove its desirable properties. Sturmfels et al. (2020) visually investigate choices of baselines and note that while one can

compare baselines through visual inspection, this approach is difficult to scale to enable a large-scale study of interpretability. Furthermore, Mohseni et al. (2021) identify user biases in the subjective rating of model saliency explanations.

To assess and compare the attributions from local explanation methods, a straightforward way is to ablate the top K features ranked by the attributions and observe the decrease of the predicted output score (Sturmfels et al., 2020). To avoid simply removing the top K features without considering the correlations among features, we can ablate the center of mass of the input instead (Ghorbani et al., 2019). Also, to avoid issues of model extrapolation on ablated inputs, we can retrain the model on the ablated data and measure the performance degradation (Hooker et al., 2018). Furthermore, local explanation methods can also be assessed by comparing their behavior between a randomly parameterized model and a trained model (Adebayo et al., 2018). Besides ablation, we can measure the quality of explanations with metrics such as (in)fidelity and sensitivity under perturbations (Yeh et al., 2019), or impact score that measures the feature importance on decision marking process (Lin et al., 2019). If the apriori of feature importance of the dataset is known, the feature importance of input across different models can also be assessed (Yang & Kim, 2019).

Although topological data analysis (TDA) is beginning to be considered for explainability purposes, the field is still nascent. One example of TDA for explainability purposes is shown by Elhamdadi et al. (2021), who use TDA to study face poses used in affective computing and find that their topology-based approach captures known patterns. van Veen (2020) proposed visual constraints on TDA output to aid interpretability, specifically for viewing global, cluster and local explanations. However, there is little work on applying TDA to the outputs of explanation methods in order to compare heterogeneous explanations or quantitatively measure consensus among explanations.

3 METHODS

3.1 TOPOLOGICAL BACKGROUND

Reeb Graph and Mapper. Consider a scalar function $f : \mathbb{M} \rightarrow \mathbb{R}$, that maps from a manifold \mathbb{M} to \mathbb{R} . The *level set* $f^{-1}(a)$ at a given scalar value a is the set of all points that have the function value a . The *Reeb graph* (Reeb, 1946) of f is computed by contracting each of the connected components of the level sets of f to a single point, resulting in a skeleton-like representation of the input.

However, many real-world data sets are available as functions defined on a set of discrete high-dimensional points rather than as continuous functions. The *Mapper* algorithm (Singh et al., 2007) computes an approximation of Reeb graph of some user-defined function (often called *lens* or *filter* function) of such data. Essentially, the Mapper algorithm divides the function range into a set of overlapping intervals and approximates the level sets to be the set of points that fall within each interval. The connected components of these approximate level sets are then computed by clustering the points that are part of a given interval. Each cluster then forms a node of the approximate Reeb graph and an edge is present between two nodes if they share one or more input points.

Persistence. Given a scalar value a , the sublevel set $f^{-1}((-\infty, a])$ is defined as the set of all points on the domain that have a function value less than or equal to a . Consider a filtration of the input that sweeps the input scalar function f with increasing function values. As the function value increases, the topology of the sublevel sets changes at the *critical points* of the function (where its gradient is zero), and remains constant at other points.

In particular, at a critical point, either a new topology is created, or some topology is destroyed. Here, topology is quantified by a class of k -dimensional cycles (or k -cycles). For example, a 0-cycle represents a connected component, a 1-cycle is a loop that represents a tunnel, and a 2-cycle bounds a void. A critical point is a creator if a new topology appears and a destroyer otherwise. Given a set of critical points c_1, c_2, \dots, c_m , one can pair up each creator c_i uniquely with a destroyer c_j which destroys the topology created at c_i . We say that a topological feature is born at c_i and it dies at c_j . The *topological persistence* (Edelsbrunner et al., 2002) of this topological feature that is created at c_i is defined as $f(c_j) - f(c_i)$, which intuitively indicates the lifetime of this feature in this sweep.

Since we work with functions defined over discrete points, we use the graph computed by the Mapper algorithm to compute the topological persistence (Edelsbrunner et al., 2002) of the features of the input. Here, the filtration is defined on the nodes and edges of the graph as follows. Each node is

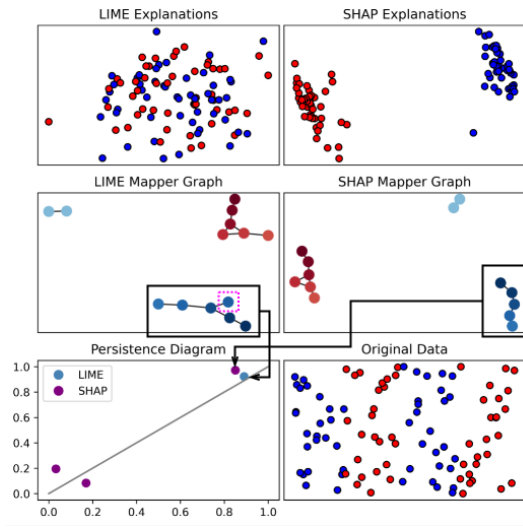


Figure 2: Comparing two sets of explanations using GALE. First, the explanation spaces are summarized using Mapper. From the Mapper output, we create persistence diagrams. Topological features correspond to points in the persistence diagram. In this example, we see that while the two Mapper outputs are similar, the Mapper output from the SHAP explanations is different from that of the LIME explanations (nodes in upper right differ, nodes in lower left overlap).

assigned a function value equal to the mean of the function values of the clustered points represented by that node. The order of the nodes added during the filtration is defined by the function value of the nodes. An edge is added during the step of the filtration as soon as both its endpoint nodes are added.

Persistence Diagrams. A persistence diagram plots the topological features as a 2-dimensional scatter plot (Edelsbrunner et al., 2002). Each point in the plot corresponds to a single feature and has x and y coordinates equal to its birth and death values respectively obtained from the extended filtration (Figure 2, lower left). Persistence diagrams provide a useful mechanism to assess the structure of scalar functions. Moreover, it has also been shown that persistence diagrams are robust to noise (Cohen-Steiner et al., 2007). We can also calculate the distance between two persistence diagrams. In this work, we use the *bottleneck distance*. Bottleneck distance measures the similarity between two persistence diagrams by finding the shortest distance b for which there exists a perfect matching between the two persistence diagrams, along with all points on the diagonal, such that any pair of matched points have at most a distance of b between them. Formally,

$$d_b(P_A, P_B) = \inf_{\gamma} \sup_x \|x - \gamma(x)\|_{\infty} \quad (1)$$

where P_A and P_B are persistence diagrams, which are multisets, $x \in P_A$, $y \in P_B$ and γ ranges over all bijections from P_A to P_B .

3.2 GALE: GLOBALLY ASSESSING LOCAL EXPLANATIONS

In this section, we introduce our approach to **globally assessing local explanations (GALE)**. Let $\mathcal{X} \subset \mathbb{R}^n$ be a data set and $P : \mathcal{X} \rightarrow [0, 1]$ a binary classification model, where $P(x) \in [0, 1]$ is the probability of $x \in \mathcal{X}$ belonging to the “1” class. Given \mathcal{X} and the model P , a local explanation method can be seen as a mapping from the data set to the explanation space $E : \mathcal{X} \rightarrow \mathbb{R}^d$, where $E(x)$ provides the “importance” of the different attributes for the classification of $x \in \mathcal{X}$. While $d = n$ for most approaches, d can be greater than n if the explanation method outputs more than a scalar value for each attribute of the input. The mapping $\mathcal{X}' = E(\mathcal{X})$ gives rise to a point set in \mathbb{R}^d , which is the explanation manifold. We can define a function $f : \mathcal{X}' \rightarrow [0, 1]$ as $f(x') = P(x)$, where $x' = E(x)$, $x \in \mathcal{X}$. In GALE, we use the function f as the lens function from which Mapper builds a summary representation G . Intuitively, f captures the relationship between the explanations and the classification probabilities. From G , we can produce its persistence diagram D .

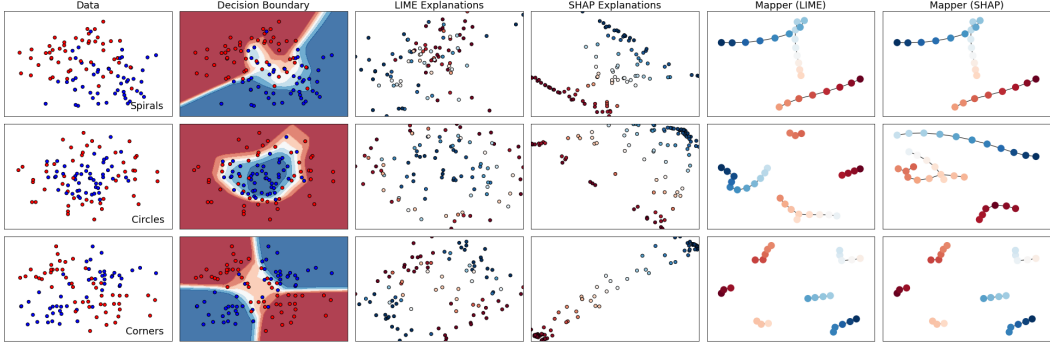


Figure 3: End-to-end pipeline showing the initial data, learned decision boundary, LIME and SHAP explanations, and the topological representations for three different 2-dimensional synthetic classification problems.

Since the explanation manifold can drastically vary across methods and parameters, a direct comparison of the geometry of these manifolds is not possible. However, studying the topology of such functions allows us to analyze how explanation methods differ and thus, we can make geometry agnostic comparisons. To do so, we can take the distance between two persistence diagrams arising from two different explanation methods. If the distance is low, then the topologies of the explanation spaces are similar.

Furthermore, each topological feature (a point in the persistence diagram) can be easily mapped back to a set of input data points in \mathcal{X} , thus allowing us to also compare how the explanation space is spread across the input data. We show an example of GALE on three synthetic data sets in Figure 3.

3.3 TUNING MAPPER PARAMETERS

Mapper requires three parameters: 1) the *resolution* r of the lens function that defines the number of intervals into which the scalar function range is divided; 2) the *gain* g , which defines the percentage overlap between successive intervals; and 3) the *clustering* algorithm used (which may carry its own parameters). The value of these parameters determines the structure of the resultant graph, and hence the persistence diagram.

In an ideal scenario, increasing the resolution and decreasing the overlap would result in the graph computed using the Mapper algorithm converging to the Reeb graph. However, in real-world data, too high a resolution or too low a gain can result in the graph being a set of disconnected nodes. In our implementation, we use agglomerative clustering. Due to the varying spaces arising from different explanation methods, we set the agglomerative clustering distance threshold parameter to a fraction of the range of the explanation space.

Depending on the data, small changes in the Mapper parameters can drastically change the resulting graph, and hence the persistence diagram. Poor selections of Mapper parameters can produce graphs which are disjointed or unstable under small perturbations to the input. Blaser and Aupetit suggest that direct measures for measuring the quality of Mapper output could be a combination of cluster quality and their consistency (Blaser & Aupetit, 2020). We are interested in identifying a set of Mapper parameters that not only produces a “stable” graph computation but also a graph with clear structure (i.e., a small number of connected components.)

To measure the aforementioned stability, we use bootstrapping (Chazal et al., 2017) to compute the confidence intervals for the bottleneck distances between persistence diagrams. Consider a set of input explanation values $E = \{e_1, e_2, \dots, e_n\}$ and the persistence diagram D generated by its Mapper graph G . G is generated from some parameter set containing the resolution, gain and clustering distance threshold. Then, for each iteration in the bootstrap, we sample with replacement from E to construct $E^* = \{e_1^*, e_2^*, \dots, e_n^*\}$, compute Mapper graph G^* and persistence diagram D^* from G^* , and calculate $d_b(D, D^*)$. Using the distribution of distances created by these iterations, we can then find the value \hat{b}_α such that

$$P(d_b(D, D^*) \geq \hat{b}_\alpha) = \alpha \quad (2)$$

A parameter set which is stable will produce a low value for \hat{b}_α . Additionally, for each iteration, we count the number of connected components in G^* . We can find the value \hat{c}_α such that

$$P(\text{connected components of } G^* \geq \hat{c}_\alpha) = \alpha \quad (3)$$

To tune the Mapper parameters, we perform a grid search over the resolution, gain and distance threshold parameters. For each parameter combination, we use 100 bootstrap iterations to estimate \hat{b}_α and \hat{c}_α for $\alpha = 0.05$. Using the estimated stability and connected components for every combination of parameters, we then iterate through the parameter set and greedily select the best combination. While this approach does not guarantee an optimal solution (e.g., we could arrive at a parameter set which has a low \hat{b}_α but a high \hat{c}_α), we find that a greedy strategy works well on a wide array of synthetic and real world data.

4 EVALUATION

4.1 COMPARING BASELINES IN GRADIENT-BASED METHODS

Here, we illustrate how GALE can be used to understand and compare different explanation methods. In particular, we aim to address a common challenge in using local explanation techniques—what should be the baseline for gradient-based explanation methods?

Baselines act as references to compare the relative importance of features for an input so that attributions can be calculated. We apply three explanation methods—Integrated Gradients (Sundararajan et al., 2017), DeepLIFT (Shrikumar et al., 2017), and SHAP with five different baselines: (1) zero baseline (an input with all values being zeros), (2) maximum distance baseline (Sturmfels et al., 2020), (3) Gaussian baseline (Smilkov et al., 2017; Sturmfels et al., 2020), (4) uniform baseline (Sturmfels et al., 2020), and (5) a trained baseline (Izzo et al., 2020), resulting in 15 explanations for each experiment.

Our goal is to illustrate how GALE can be used to identify explanation methods that behave differently. To do so, we generate a synthetic data set with 5 features and determine the labels with an extremely simple logic—the input will be labeled as “1” if any of the columns contain a “zero” as value. Otherwise, they are labeled as “0”. When working with real data sets, it is normal for them to contain zero values that may provide important information. Under our construction, only zero values are important to the classification. Thus we expect the zero baseline is the only reference that is not a neutral input to the classifier. We created 20 different synthetic data sets under the aforementioned conditions, where each set varied in its rate of “0” labeled instances.

We first calculate a pairwise distance matrix for each of the 20 data sets. Then, we average the distance matrices, and show the result in Figure 4. We can see that there are three rows and columns which have much greater average distances (yellow color) than others, meaning that the persistence diagrams generated by these explanations differ greatly as compared to the other explanations. Such observation is consistent with the fact that zero baselines produce different feature attributions by treating the important values as neutral references. Furthermore, when we explore the corresponding topological skeletons, we observe that these explanations typically have a more scattered graph compared with the others. Such visual evidence provide a sense that the explanation values differ greatly among the inputs in the data set. This evidence can be found statistically as well. When computing the variances of the explanation values among different outcomes, the variances are generally high among the outputs with scattered topological graphs.

4.2 EVALUATING TOPOLOGICAL STABILITY

Many methods, like LIME, rely on stochastic behavior, which can produce different explanations, even when given the same parameters. It is important that GALE is resistant to sampling variability incurred by local explanation methods. In this section, we benchmark GALE’s stability with regards to the variation induced by local explainability methods. Additionally, we show how our greedy approach to Mapper parameter selection aids in finding stable topological representations.

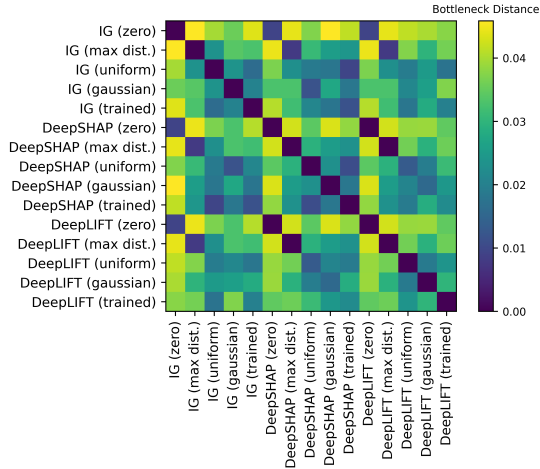


Figure 4: Pairwise distance matrix showing the average bottleneck differences among gradient-based explanation methods across 20 synthetic datasets. We see that the methods which use the zero-baseline (rows/columns denoted with “zero”) have higher bottleneck distances to the explanations using other baselines.

Data	Row-Wise Distance		Connected Components	
	<i>Greedy</i>	<i>Fixed</i>	<i>Greedy</i>	<i>Fixed</i>
Spirals	0.00	0.83	2.0	4.1
Circles	1.88	2.56	6.1	10.5
Corners	1.22	1.61	9.0	14.6
Lin. Sep.	0.00	0.00	4.0	4.0
Toy	0.03	0.00	2.0	3.0
Toy-Flip	0.00	0.00	2.0	2.2
Toy-Int.	0.16	0.00	2.0	5.0
Diabetes	0.38	2.70	4.3	36.4

Table 1: Average row-wise sum of distance matrix and average connected component counts after using parameters found by our greedy parameter search. We find that our parameter tuning procedure produces more connected and stable Mapper output.

To determine the robustness of GALE against local explanation method variability, for each synthetic data set and the diabetes data set, we run LIME 30 times. We restrict our experiments to these data sets for computational purposes. For each run, we generate two Mapper outputs: one which uses the Mapper parameters found via a greedy search, as described in Section 3.2, and one which uses fixed parameters: a resolution of 15, a gain of 0.3 and an agglomerative distance threshold of 0.3 times the range of the explanation space. We then generate two bottleneck distance matrices – one for distances between the optimized Mappers and another for distances between the fixed parameter Mappers. In Table 1, we report the average row-sum for each of these matrices for each synthetic data set, as well as the average number of connected components in the Mapper graphs. The higher the row-sum, the less consensus there is among the explanation topologies produced by Mapper. We see that by using our greedy Mapper parameter search, our output is resistant to variation in LIME output.

4.3 USING GALE TO TUNE EXPLAINABILITY METHOD PARAMETERS

Many local explainability methods use a variety of parameters which guide the output. For example, in SHAP, one can specify parameters such as the number of times to re-evaluate the model when explaining each prediction in Kernel SHAP, or the limit on the number of trees used in Tree SHAP. Choosing an appropriate set of parameters is important, as there can be significant variability in explanations coming from even just slightly different parameters (Visani et al., 2020). Here, we show

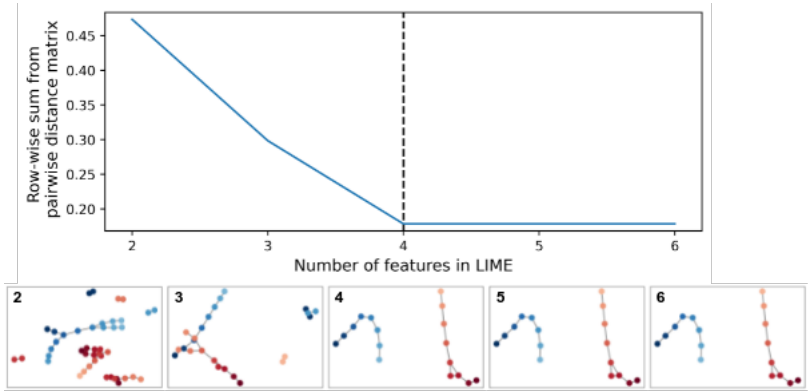


Figure 5: Row-wise sum of distance matrix for Mapper graphs computed on LIME explanations of different feature counts. Since the average row-wise distance stays constant when using 4 or more features for LIME, we know that the produced topologies are similar for these parameter values.

how to use GALE to tune the parameters of a local explainability method – LIME. We consider a toy synthetic dataset with six independent features, where only four are used for the prediction. We train a random forest classifier on this data using the default sklearn parameters. Our goal is to determine the number of features to use for our LIME explanation, which we know to be four. We set the size of the neighborhood to 50 observations.

Using our Mapper parameter tuning technique to find appropriate parameters, we compute the corresponding graph for explanations using k features, where $k \in \{2, \dots, 6\}$. Then, we calculate the bottleneck distance between each of the persistence diagrams. Next, we calculate the row-wise sum of the distance matrix. Persistence diagrams which differ strongly from others will have a high row-wise sum. We plot the row-wise sum in Figure 5. We see that when the number of features is four or higher, the returned topological signatures are identical. Thus, one may conclude that the addition of a fifth or sixth feature in the LIME explanations has no effect on the overall topology of the explanations. Likewise, by only including two or three features in LIME, one may not be estimating appropriate explanations.

5 CONCLUSION

We present GALE, a topology-based framework to globally summarize the outputs of local explanation methods. To do so, we compute a topological skeleton that captures the relationship between the explanations and model predictions. Then, using this skeleton, we find its persistence diagram, which represents the topological features in the explanation space. Finally, we calculate the distance between the different explanations’ corresponding persistence diagrams to derive a measure of similarity between explanations. GALE (1) allows for easy comparison of heterogeneous local explanations, (2) is resistant to outliers and variation, and (3) can be used to optimize explanation method parameters. While GALE sheds no light on what are the *correct* explanations for a given prediction problem, it is a powerful and simple tool to quantitatively compare heterogeneous local explanations, compared to visual inspection or user studies.

Limitations. Although topological data analysis shows promise in tackling challenges encountered in XAI applications, there are still several shortcomings which we plan to address in the future. While we focus on binary classification with tabular data, which is a canonical machine learning task, extending GALE to consider multiclass problems, as well as text or image models, will be a key line of future work. Additionally, although we did not evaluate GALE on regression problems, the approach would be similar, except for changing the lens function from the class predictions to the regression output.

ACKNOWLEDGMENTS

We would like to thank Bruno Coelho and Richen Du for their support in running preliminary experiments. This collaboration has been supported by a grant from Capital One. Silva’s research has also been supported by NASA; NSF awards CNS-1229185, CCF1533564, CNS-1544753, CNS-1730396, CNS-1828576, CNS-1626098; and DARPA. Nonato’s research is supported by Fapesp-2013/07375-0 and CNPq-307184/2021-8. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of DARPA, NSF, NASA, Fapesp, CNPq, or Capital One.

REFERENCES

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *arXiv preprint arXiv:1810.03292*, 2018.
- Brian Barr, Ke Xu, Claudio Silva, Enrico Bertini, Robert Reilly, C Bayan Bruss, and Jason D Wittenbach. Towards ground truth explainability on tabular data. *arXiv preprint arXiv:2007.10532*, 2020.
- Nello Blaser and Michael Aupetit. Research directions to validate topological models of multi-dimensional data. In *NeurIPS 2020 Workshop on Topological Data Analysis and Beyond*, 2020.
- Frédéric Chazal, Brittany Fasy, Fabrizio Lecci, Bertrand Michel, Alessandro Rinaldo, Alessandro Rinaldo, and Larry Wasserman. Robust topological inference: Distance to a measure and kernel distance. *The Journal of Machine Learning Research*, 18(1):5845–5884, 2017.
- David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams. *Discrete Comput. Geom.*, 37:103–120, 2007.
- Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.
- Herbert Edelsbrunner, David Letscher, and Afra Zomorodian. Topological Persistence and Simplification. *Discrete Comput. Geom.*, 28(4):511–533, 2002.
- Hamza Elhamedi, Shaun Canavan, and Paul Rosen. Affectivetda: Using topological data analysis to improve analysis and explainability in affective computing. *IEEE Transactions on Visualization and Computer Graphics*, 2021.
- Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3681–3688, 2019.
- Johannes Haug, Stefan Zürn, Peter El-Jiz, and Gjergji Kasneci. On baselines for local feature attributions. *arXiv preprint arXiv:2101.00905*, 2021.
- Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. *arXiv preprint arXiv:1806.10758*, 2018.
- Cosimo Izzo, Aldo Lipani, Ramin Okhrati, and Francesca Medda. A baseline for Shapely values in MLPs: from missingness to neutrality. *arXiv preprint arXiv:2006.04896*, 2020.
- Jeya Vikranth Jeyakumar, Joseph Noor, Yu-Hsi Cheng, Luis Garcia, and Mani Srivastava. How can I explain this to you? an empirical study of deep neural network explanation methods. *Advances in Neural Information Processing Systems*, 33:4211–4222, 2020.
- Zhong Qiu Lin, Mohammad Javad Shafiee, Stanislav Bochkarev, Michael St Jules, Xiao Yu Wang, and Alexander Wong. Do explanations reflect decisions? A machine-centric strategy to quantify the performance of explainability algorithms. *arXiv preprint arXiv:1910.07387*, 2019.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 4765–4774. Curran Associates, Inc., 2017.

- Sina Mohseni, Jeremy E. Block, and Eric D. Ragan. Quantitative evaluation of machine learning explanations: A human-grounded benchmark. In Tracy Hammond, Katrien Verbert, Dennis Parra, Bart P. Knijnenburg, John O'Donovan, and Paul Teale (eds.), *IUI '21: 26th International Conference on Intelligent User Interfaces, College Station, TX, USA, April 13-17, 2021*, pp. 22–31. ACM, 2021. doi: 10.1145/3397481.3450689.
- Georges Reeb. Sur les points singuliers d'une forme de Pfaff complètement intégrable ou d'une fonction numérique. *Comptes Rendus des séances de l'Académie des Sciences*, 222(847-849):76, 1946.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp. 1135–1144, 2016.
- Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3145–3153. JMLR. org, 2017.
- Gurjeet Singh, Facundo Mémoli, and Gunnar E Carlsson. Topological methods for the analysis of high dimensional data sets and 3d object recognition. In *Eurographics Symp. on Point-Based Graphics*, volume 91, pp. 100, 2007.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Kacper Sokol and Peter Flach. Explainability fact sheets: a framework for systematic assessment of explainable approaches. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 56–67, 2020.
- Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Visualizing the impact of feature attribution baselines. *Distill*, 5(1):e22, 2020.
- Mukund Sundararajan and Amir Najmi. The many shapley values for model explanation. In *International Conference on Machine Learning*, pp. 9269–9278. PMLR, 2020.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pp. 3319–3328. PMLR, 2017.
- Hendrik Jacob van Veen. Novel topological shapes of model interpretability. In *NeurIPS 2020 Workshop on Topological Data Analysis and Beyond*, 2020.
- Giorgio Visani, Enrico Bagli, Federico Chesani, Alessandro Poluzzi, and Davide Capuzzo. Statistical stability indices for lime: obtaining reliable explanations for machine learning models. *Journal of the Operational Research Society*, pp. 1–11, 2020.
- Daniel S. Weld and Gagan Bansal. The challenge of crafting intelligible intelligence. *Commun. ACM*, 62(6):70–79, 2019.
- Mengjiao Yang and Been Kim. Benchmarking attribution methods with relative feature importance. *arXiv preprint arXiv:1907.09701*, 2019.
- Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Sai Suggala, David I Inouye, and Pradeep Ravikumar. On the (in) fidelity and sensitivity for explanations. *arXiv preprint arXiv:1901.09392*, 2019.