
gLIME: A NEW GRAPHICAL METHODOLOGY FOR INTERPRETABLE MODEL-AGNOSTIC EXPLANATIONS

A PREPRINT

✉ **Zoumpolia Dikopoulou***
AIDEAS OÜ,
Tallinn, 10117, Estonia.
dikopoulia@gmail.com

✉ **Serafeim Moustakidis**
AIDEAS OÜ,
Tallinn, 10117, Estonia.
s.moustakidis@aideas.eu

Patrik Karlsson
AIDEAS OÜ,
Tallinn, 10117, Estonia.
p.karlsson@aideas.eu

July 22, 2021

ABSTRACT

Explainable artificial intelligence (XAI) is an emerging new domain in which a set of processes and tools allow humans to better comprehend the decisions generated by black box models. However, most of the available XAI tools are often limited to simple explanations mainly quantifying the impact of individual features to the models' output. Therefore, human users are not able to understand how the features are related to each other to make predictions, whereas the inner workings of the trained models remain hidden. This paper contributes to the development of a novel graphical explainability tool that not only indicates the significant features of the model, but also reveals the conditional relationships between features and the inference capturing both the direct and indirect impact of features to the models' decision. The proposed XAI methodology, termed as *gLIME*, provides graphical model-agnostic explanations either at the global (for the entire dataset) or the local scale (for specific data points). It relies on a combination of local interpretable model-agnostic explanations (LIME) with graphical least absolute shrinkage and selection operator (GLASSO) producing undirected Gaussian graphical models. Regularization is adopted to shrink small partial correlation coefficients to zero providing sparser and more interpretable graphical explanations. Two well-known classification datasets (BIOPSY and OAI) were selected to confirm the superiority of *gLIME* over LIME in terms of both robustness and consistency/sensitivity over multiple permutations. Specifically, *gLIME* accomplished increased stability over the two datasets with respect to features' importance (76%-96% compared to 52%-77% using LIME). *gLIME* demonstrates a unique potential to extend the functionality of the current state-of-the-art in XAI by providing informative graphically given explanations that could unlock black boxes.

Keywords AI explainability · interpretability · model-agnostic explanations · graph models

1 Introduction

Superhuman capacity has been recently demonstrated by artificial intelligence (AI) leading to a widespread adoption of AI systems in multiple domains including healthcare, industry 4.0 and finance. However, this improved predictive performance typically comes with increased model complexity LeCun et al. [2015]. The great majority of recent powerful machine learning algorithms are 'black box' approaches with the rationale behind their decision-making mechanism being hard to understand and interpret. The ambiguity with respect to the models' inner workings and the fact that their decisions cannot be interpreted make these systems difficult to be trusted by the end-users, especially in critical domains such as in healthcare. Meeting the need for trustworthy, fair and robust AI decisions, explainable AI (XAI) Gunning and Aha [2019] has emerged as a new scientific field that focuses on the understanding and interpretation of AI systems' behavior.

*<https://www.aideas.eu/>

So far, XAI has been approached from different view-points with respect to the type of data (tabular, images or text) or the interpretability scale Linardatos et al. [2021]. As far as the interpretability scale, some techniques provide explanations for individual instances (local scale), whereas there are XAI tools or libraries used to explain the behavior of the model on the whole dataset (global scale). Another important classification of XAI tools is related to the type of algorithms/networks in which the explainability analysis is applied. There are XAI techniques that are model specific (restricted to a specific machine learning model or to specific family of models) and other that are model-agnostic capable to be applied to any model.

Post-hoc explainability refers to a specific category of XAI that encompasses techniques that explain the decisions of already trained black-box models. A considerable amount of experiments and scientific work has been devoted on the explainability of deep learning models and thus a variety of model-specific XAI tools has been proposed including DeepLIFT Shrikumar et al. [2017], Class Activation Maps (CAMs), first introduced in Zhou et al. [2016], and Grad-CAM Selvaraju et al. [2017]. Among the post-hoc model-agnostic techniques, the local interpretable model-agnostic explanations (LIME) method Ribeiro et al. [2016] is one of the most popular methods for black-box models that generates interpretations at the local scale (for single instances). LIME is a simple but powerful technique that derives explanations utilizing simulated randomly-sampled data around the neighbourhood of an input instance. However, LIME has been proved sensitive to these randomly generated permutations leading to unstable interpretations Garreau and Luxburg [2020]. Shapley Additive explanations (SHAP) Lundberg and Lee [2017] is another well-known game-theory inspired technique that estimates the importance of each feature on individual predictions, demonstrating both accuracy and consistency. Overall, all the aforementioned model-agnostic techniques, including SHAP, do not take feature dependence into account and in some cases produce non-intuitive feature importance values.

Current post-hoc model-agnostic XAI techniques are limited to a very specific view-point of XAI where feature importance values are calculated and visualized with bar graphs or other similar visualization tools. To the best of our knowledge, none of the available techniques is capable of identifying the relationships between features and the possible indirect effect of features to the models' output. This paper contributes to the development of a novel graphical explainability tool that not only indicates the significant features of the model, but also reveals the conditional relationships among features capturing and visualizing both the direct and indirect impact of features to the models' decision. The proposed XAI methodology, termed as gLIME, visualizes model-agnostic explanations with intuitive interpretable graphs at either the global (explanations for the entire dataset) or the local scale (explanations for specific data points). It relies on a combination of local interpretable model-agnostic explanations (LIME) with graphical least absolute shrinkage and selection operator (GLASSO) Epskamp and Fried [2018], Meinshausen et al. [2006] producing undirected Gaussian graphical models. In this paper we demonstrate the effectiveness of gLIME at the local scale and we compare it with LIME that shares similar characteristics. An extensive experimental analysis has been performed using two well-known classification datasets to confirm the superiority of gLIME over LIME in terms of both robustness and consistency/sensitivity over multiple permutations.

This paper is organized as follows. Section 2 gives an overview of the proposed gLIME methodology presenting its main characteristics and features. Results on multiple experiments are provided in Section 3, whereas conclusions are drawn in Section 4.

2 Methods

The fundamental idea behind the gLIME algorithm is focused on the local model-agnostic explainability and interpretability increasing the trust and the stability of the produced model. This signifies that the proposed methodology attempts to explain and interpret individual predictions of a model-agnostic method, in which gLIME can be applied to any supervised regression or classification model. This algorithm consists of two parts; the first one prepares the data by generating additional data around the selected observation and the second part, estimates a graph model. Generally, a graph or a network G is an abstract model which represents complex phenomena, and it consists of two components, nodes and edges Dorogovtsev and Mendes [2013]. Nodes or vertices (V) represent entities (features) and they are visualized as circles; while, edges or links (E) connect the nodes between them representing their relationships. When graph models are estimated from data structures, the edges may represent unknown statistical relationships such as: correlations, covariances, partial correlations, regression coefficients, factor loadings, etc. Hevey [2018].

Therefore, gLIME incorporates graphs to explain and interpret the decisions of trained models for three main reasons: i) to identify the relevant features that affect the model's prediction, ii) highlight hidden and/or important relationships among features quantifying the strength of the relationships and iii) determine significant path routes presenting how the information flows from one node to the end node (the predicted feature) traversing the intermediate significant nodes that can affect indirectly the predicted output.

Below the steps of *gLIME* are described. Before dividing the dataset into training and testing, the dataset must be clean from missing values. Then, a selected ML model is applied to the training data for a classification or regression problem. In our experiment, Non linear Support Vector Machines (SVMs) were utilized to perform the classification tasks. From the testing data, an observation is chosen and it is permuted m times to create replicated feature data with minor value variations. Next, a similarity distance measure is applied to calculate the distance between the initial observation and the permuted observations. Particularly, the distance measure is converted to a similarity value with the use of an exponential kernel which by default is adjusted to 0.75 times the square root of the number of features. Afterwards, the ML model is applied to these new points to predict the outcomes (scores or probabilities to a certain class). Consequently, a new dataset is created which includes the permuted data points of the features that are close to the original observation and the corresponding predictions (scores).

Afterwards, the network model is estimated to explain and interpret the influential features of the predicted outcome, the important interconnections among the features and the significant path route of the model. The proposed algorithm introduces the undirected graphical models as the most well-known frameworks for constructing a network model revealing the straightforward relationships between observable features. Particularly, when data follow a multivariate normal distribution, the produced model is called Gaussian graphical model (GGM; Lauritzen [1996]) and belongs to a generalized class of statistical models named pairwise Markov random fields (PMRF; Koller and Friedman [2009]). The produced GGM estimates can be standardized, visualized and easier interpreted as partial correlation coefficients Borsboom and Cramer [2013], McNally [2016]. In particular, partial correlation coefficients fluctuate between -1 and 1 and reveal the remaining linear dependency among two variables, after conditioning on all other variables in the dataset Epskamp and Fried [2018].

Partial correlations can be directly computed from the inverse of a variance–covariance matrix in which each element represents a weight–edge indicating the strength of connection between two nodes Epskamp and Fried [2018]. Let’s assume that $Y^T = [Y_1 Y_2 \dots Y_p]$ represents the response vector of a random subject of p features. Supposing that y vector is centred and follows a multivariate normal density with some $p \times p$ variance-covariance matrix Σ , $Y \sim N_p(0, \Sigma)$. The partial correlation coefficients are estimated by calculating the inverse of Σ (known also as a precision matrix), $K = \Sigma^{-1}$. The element k_{ij} can be standardized to derive the partial correlation coefficient among variables Y_i and Y_j after conditioning on all other variables in Y , $Y_{-(i,j)}$, Lauritzen [1996], $Cor(Y_i, Y_j | Y_{-(i,j)}) = \frac{-k_{ij}}{(\sqrt{k_{ii}})\sqrt{k_{jj}}}$.

If the partial correlation coefficient is exactly zero, this signifies a conditional independence between two variables after controlling for all other variables in the model and therefore, no edge is drawn between these two nodes.

However, in practice, if two features are conditionally independent, small partial correlations (close to zero) are estimated which are called spurious or false positives Costantini et al. [2015]. In order to control the spurious connections in the precision matrix, a statistical regularization technique originating in the field of machine learning is applied. This technique is known as ‘least absolute shrinkage and selection operator’ (lasso; Tibshirani [1996]) which shrinks small partial correlation coefficients exactly to zero. The graphical lasso is a regularization framework for estimating the covariance matrix Σ under the assumption that precision matrix K is sparse Meinshausen et al. [2006]. The graphical lasso problem maximizes the penalized (l_1 -regularized) log-likelihood:

$$\text{maximize}_{K>0} f(K) := \log \det(K) - (SK) - \lambda \|K\|_1 \tag{1}$$

where S denotes the sample covariance matrix, λ (lambda) is a nonnegative tuning parameter controlling the amount of l_1 shrinkage and $\|K\|_1$ is the l_1 norm (the sum of the absolute values of the elements of Σ^{-1}). Since λ regulates the sparsity of the network, various values of λ provide different network structures Zhao and Yu [2006] indicating that a group of networks ranging from a fully connected network (λ_{min}) to an empty network (λ_{max}) are estimated. Typically, a logarithmically spaced range of tuning parameters in which $\lambda_{min} = R\lambda_{min}$ where $R = 0.01$ by default. Subsequently, the network that minimizes the Extended Bayesian Information Criterion (EBIC; Drton and Perlman [2004]) is characterized as optimal network signifying that fits better into the data. As shown in (2), the EBIC adds an extra penalty, the hyperparameter γ (gamma) to control the model complexity Foygel and Drton [2010] and it is set manually from 0 to 0.5, L indicates the log-likelihood, n the sample size, E the number of non-zero edges and p the number of nodes.

$$EBIC = -2L + E \log n + 4\gamma E \log p \tag{2}$$

It is conducted that the combination of lasso regularization with EBIC model selection provides the true network structure, particularly when a sparser network is estimated (Epskamp and Fried, 2018; Foygel and Drton, 2010). A $p \times p$ undirected weighted matrix is returned in which the conditional dependent relationships among variables are presented and stronger connections among features are identified. To explain the features’ importance on the predicted

outcome, we rank the features that are directly connected to the predicted feature. Finally, the path that traverses the most important features is computed from the end node (the predicted feature). Specifically, it searches which node is highly connected to the end node, then the selected node searches at its neighborhood the node with the strongest connection. The algorithm stops when i) all nodes of the model are traversed or ii) the edge-weight among node i and node j (w_{ij}) is lower than a specific value (by default, it stops when $w_{ij} < 0.1$). For the sake of simplicity, the steps of gLIME algorithm are summarized in Pseudocode I.

Table 1: The gLIME pseudocode

Pseudocode I: <i>The gLIME algorithm</i>
1. Select an observation (row) from the testing data.
2. Permute the observation m times (default 5000 times).
3. Calculate the distance from all permutations to the selected observation.
4. Use a ML model to predict the outcome of all permuted observations.
5. Save the new permuted dataset with the predicted outcome.
6. Create a vector of λ values (default 100 values).
7. For each λ , apply the Equation (1), the graphical lasso methodology to create a graphical model.
8. Calculate the EBIC in Equation (2) for each estimated graphical model.
9. Select the graphical model that minimizes the EBIC and return the undirected weighted matrix, $p \times p$.
10. Model explain ability: Rank the features that are directly connected to the predicted outcome according to their weights.
11. a. Model interpretability: Determine the higher connections among two features of the model.
b. Model interpretability: Find the significant path route to present how the information flow though highly relevant features.

3 Results and Discussion of Results

In this study, the gLIME algorithm was applied in two datasets (BIOPSY and OAI) to explain which of the features influenced significantly the predicted output, understand better the estimated graph model by interpreting the strongest connections among the features and indicate the significant path in the graph presenting how the predicted output could be indirectly affected by other features of the network. Every dataset was divided into the training and testing dataset. From the testing data, four random observations were selected and for each observation, ten permuted datasets were generated including 5000 permuted observations. In total, 120 permuted datasets were produced, forty for each problem. These permutations were necessary as well to confirm the stability of the gLIME algorithm compared to the results of LIME.

Specifically, the BIOPSY dataset includes biopsies of breast tumors of 699 patients (<https://github.com/cran/MASS/blob/master/data/biopsy.rda>). Each of the nine attributes (V_1 : Clump thickness, V_2 : Uniformity of cell size, V_3 : Uniformity of cell shape, V_4 : Marginal adhesion, V_5 : Single epithelial cell size, V_6 : Bare nuclei, V_7 : Bland chromatin, V_8 : Normal nucleoli and V_9 : Mitoses) has been scored on a scale of 1 to 10, and the outcome was classified in two classes ‘benign’ or ‘malignant’. The OAI dataset (<https://nda.nih.gov/oai/>) focused on Osteoarthritis problem; specifically, the goals of the OAI are to provide resources to enable a better understanding of prevention and treatment of knee osteoarthritis. For this study, the OAI dataset consists of 3873 patients and forty features were examined; while the outcome was categorized in two categories ‘healthy’ and ‘not healthy’.

Table 2 illustrates the graphical lasso estimates of the gLIME algorithm applied on the first permuted dataset of the first observation. As it is observed, the derived matrix is a 10×10 undirected weighted matrix in which the number ten represents the overall number features (nine inputs and one output). Moreover, it was also inferred the strongest

conditional associations between two features after conditioning on all other features in the model. The indicative strongest positive conditional associations were identified between the Uniformity of cell size and Uniformity of cell shape ($w_{2,3}=0.594$), the Uniformity of cell size and the Single epithelial cell size ($w_{2,5}=0.231$) and the Single epithelial cell size and Mitoses ($w_{5,9}=0.178$). Thus, the strongest negative conditional associations were detected among input features (Bare nuclei, Clump thickness and Bland chromatin) and the predicted outcome which was classified as ‘benign’ after conditioning on all other features in the model, $w_{6,01}=-0.395$, $w_{1,01}=-0.304$ and $w_{7,01}=-0.216$. The negative conditioning connection signifies the reverse relation between two features after conditioning on all other features in the model. For instance, the $w_{6,01}=-0.395$ determines the smaller the Bare nuclei attribute is, the healthier the patient (higher possibility to be benign cancer) after conditioning on all other observed attributes and vice-versa. Moreover, six out of forty-five connections were distinguished as zero implying the conditional independence among features $w_{1,7}=w_{1,8}=w_{2,6}=w_{3,9}=w_{5,01}=w_{6,8}=0$.

Table 2: The weight adjacency matrix using *gLIME* in the first permuted dataset of the first observation applied on ten features (nine inputs: V_1 - V_9 and one output: O_1).

	V1	V2	V3	V4	V5	V6	V7	V8	V9	O1
V1	0	0.046	0.095	-0.06	0.015	0.048	0	0	0.002	-0.304
V2	0.046	0	0.594	0.107	0.231	0	0.142	0.081	0.026	-0.093
V3	0.095	0.594	0	0.001	0.081	0.057	0.012	0.123	0	-0.153
V4	-0.06	0.107	0.001	0	0.059	0.171	0.145	0.089	0.098	-0.12
V5	0.015	0.231	0.081	0.059	0	0.037	0.031	0.118	0.178	0
V6	0.048	0	0.057	0.171	0.037	0	0.098	0	-0.087	-0.395
V7	0	0.142	0.012	0.145	0.031	0.098	0	0.171	-0.066	-0.216
V8	0	0.081	0.123	0.089	0.118	0	0.171	0	0.085	-0.058
V9	0.002	0.026	0	0.098	0.178	-0.087	-0.066	0.085	0	-0.136
O1	-0.304	-0.093	-0.153	-0.12	0	-0.395	-0.216	-0.058	-0.136	0

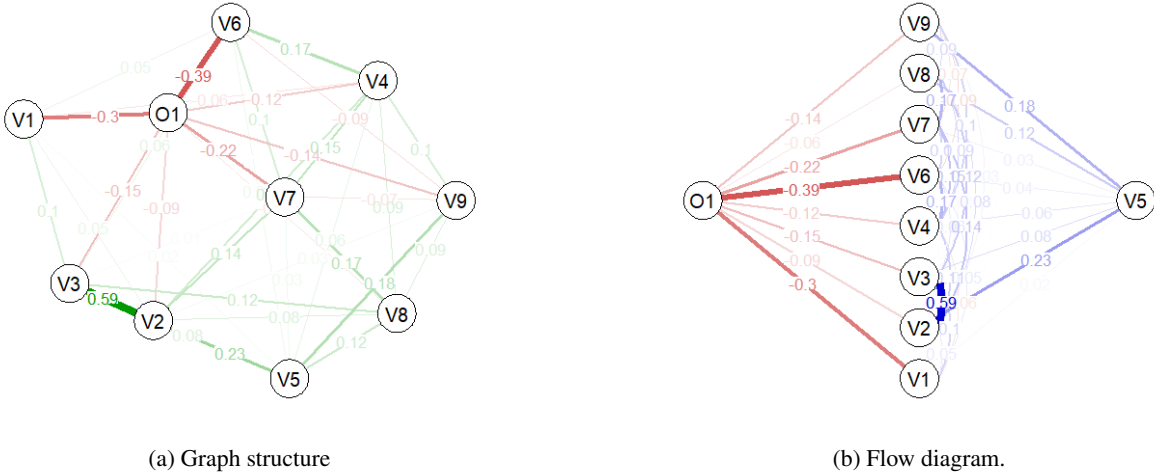


Figure 1: Visual representation of the regularized partial correlation matrix of Table 1 as a partial network structure. Circles represent the observed features and links represent the regularized partial correlation coefficients among two features after conditioning on all other features of the model. Green/blue and red edges denote positive and negative associations, respectively. Wider and more saturated edges indicate higher strengths among nodes.

Since every graph can be described by an adjacency matrix, Figure 1 illustrates the corresponded graphical model of Table 2. Every weight of the graph is colored according to its strength. Specifically, green or blue colors represent positive regularized partial associations and red characterize negative regularized partial associations. Thus, saturated edges define the magnitude of the partial strength among two nodes, i.e. lighter color display weaker regularized partial correlations and darker pigment illustrate higher regularized partial associations. There are many ways to visualize a graph structure. In Figure 1, the weighted matrix of Table 2 has been presented as a) a simple graph in which strongly connected nodes were placed closer reflecting the full picture of the model; while, weakly related nodes were set

away from each other and b) a flow diagram in which the node of interest (O1) was placed to the left, while the nodes that were directed connected to the node of interest were placed vertically in a new layer (similar to a neural network structure) and so on, presenting clearly the significant features related to the node of interest.

Table 3: The ranking positions of the features that are directly linked with the prediction outcome of ten permuted data (1A – 1J) using *gLIME* and *LIME* concerning the BIOPSY dataset.

		Rankings																			
		1A		1B		1C		1D		1E		1F		1G		1H		1I		1J	
		<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>	<i>gLIME</i>	<i>LIME</i>
Variables	V1	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
	V2	7	7	8	8	7	7	8	5	7	4	7	4	8	4	6	3	7	4	8	4
	V3	4	4	4	4	3	3	4	3	4	3	4	3	4	3	4	4	4	3	4	3
	V4	6	6	5	5	5	5	5	6	6	7	6	7	6	7	7	7	6	7	6	5
	V5	9	9	9	9	9	9	9	8	9	8	9	8	9	8	9	8	9	6	9	8
	V6	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	V7	3	3	3	3	4	4	3	4	3	5	3	5	3	5	3	5	3	5	3	6
	V8	8	8	7	7	8	8	7	9	8	9	8	9	7	9	8	9	8	9	7	9
	V9	5	5	6	6	6	6	6	7	5	6	5	6	5	6	5	6	5	8	5	7

Table 4: Biopsies of breast tumors problem. The Kendall’s tau coefficients of ten permuted data of the first observation. The τ_b of *gLIME* and *LIME* are highlighted with green and blue color, respectively.

		Kendal's tau coefficients of ten permuted data									
		<i>gLIME</i>									
		1A	1B	1C	1D	1E	1F	1G	1H	1I	1J
LIME	1A	1.000	.889**	.889**	.889**	1.000**	1.000**	.944**	.944**	1.000**	.944**
	1B	.889**	1.000	.889**	1.000**	.889**	.889**	.944**	.833**	.889**	.944**
	1C	.889**	.889**	1.000	.889**	.889**	.889**	.833**	.833**	.889**	.833**
	1D	.722**	.722**	.833**	1.000	.889**	.889**	.944**	.833**	.889**	.944**
	1E	.722**	.611*	.722**	.889**	1.000	1.000**	.944**	.944**	1.000**	.944**
	1F	.722**	.611*	.722**	.889**	1.000**	1.000	.944**	.944**	1.000**	.944**
	1G	.722**	.611*	.722**	.889**	1.000**	1.000**	1.000	.889**	.944**	1.000**
	1H	.667*	.556*	.667*	.833**	.944**	.944**	.944**	1.000	.944**	.889**
	1I	.556*	.556*	.667*	.833**	.833**	.833**	.833**	.778**	1.000	.944**
	1J	.611*	.611*	.722**	.889**	.889**	.889**	.889**	.833**	.833**	1.000

** . Correlation is significant at the 0.01 level (2-tailed).
 * . Correlation is significant at the 0.05 level (2-tailed).

The features that were directly connected with the prediction outcome were ranked according to their absolute edge weight. Table 3 includes the ranking results of the first observation using *gLIME* and *LIME* over ten permuted data (1A – 1J). Overall, these results indicated that the most important features are V6 and V1; however, after the second position, small differences were observed. To compare the stability and the consistency of *gLIME* and *LIME*, Kendall’s tau coefficient (τ_b) was applied Kendall [1948]. Table 4 illustrates Kendall’s tau coefficients to determine the relationships among the permuted data of the first observation towards *gLIME* and the relationships (τ_b) of *LIME*, respectively. The average (τ_b) correlation of *gLIME* highlighted with green color, indicated that there was a very strong, positive relationship that was statistically significant (*gLIME*: $\tau_b=.934$). On the other hand, the average Kendall’s correlation of *LIME* depicted with a blue color, determined that there was a strong, positive relationship that

was statistically significant (LIME: $\tau_b=.819$). Consequently, the average stability of the first observation using gLIME was higher 93.4% comparing to LIME (81.9%).

For the sake of completeness, Table 5 summarizes the robustness of the selected observations (1st, 6th, 100th and 600th). Specifically, forty permuted datasets were generated, ten permuted datasets per observation. In total, both algorithms presented high robustness in all observations (over 81.9%); but, the average stability of each observation in gLIME revealed that the results were less sensitive comparing to LIME. Thus, in the last column of Table 5, the average stability of gLIME has shown higher robustness coefficients over the selected observations in BIOPSY dataset in which a small number of features were tested. In addition, the corresponded summarized stability coefficients of the OAI problem applied on four random observations (1st, 10th, 100th and 200th) using forty permuted datasets (Table 6). The results revealed that gLIME performed better in terms of robustness since the average Kendall’s correlation coefficients were exceeded in all of the selected observations (69.4% - 81.3%) comparing to LIME stability results (45.5% - 58.1%). Moreover, these results highlighted the inefficiency of LIME regarding the consistency estimates when the number of features were increased.

Table 5: Summarized stability coefficients of four observations which were selected from the BIOPSY testing set using gLIME and LIME algorithm.

	Observations				Mean
	1st	6th	100th	600th	
gLIME	93.4%	95.2%	97.2%	98.1%	96.0%
LIME	81.9%	86.3%	86.0%	85.0%	84.8%

Table 6: Summarized stability coefficients of four observations which were selected from the OAI testing set using gLIME and LIME algorithm studying forty features.

	Observations				Mean
	1st	10th	100th	200th	
gLIME	79.3%	79.5%	69.4%	81.3%	75.9%
LIME	52.4%	45.5%	50.3%	58.1%	51.6%

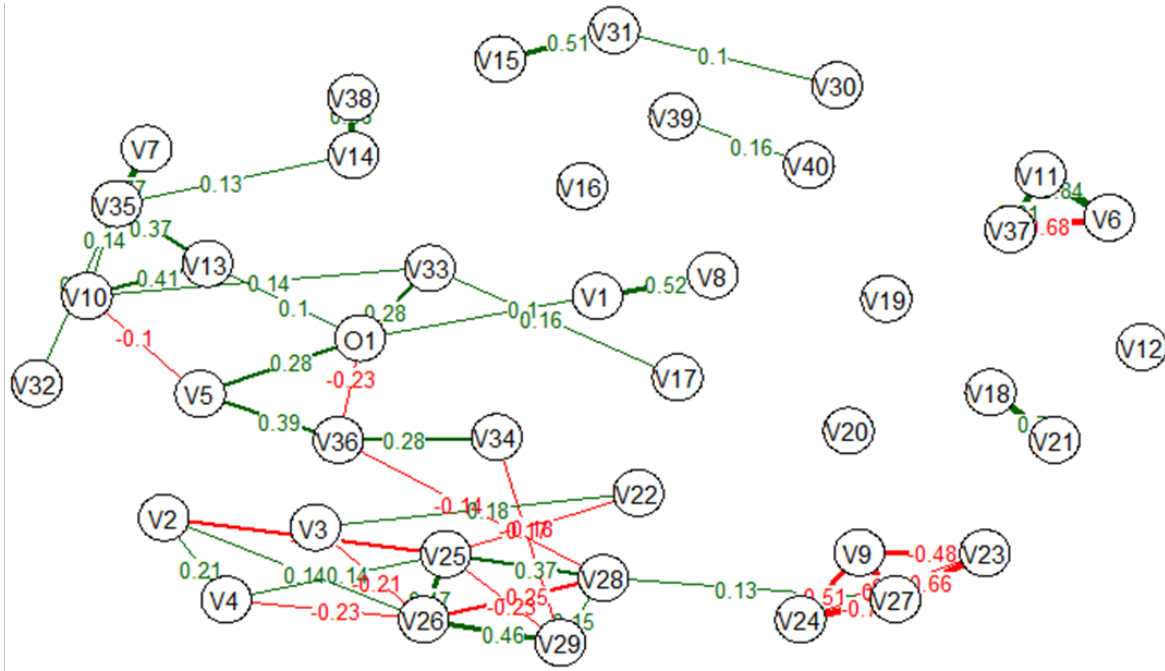
Another advantage of gLIME was the ability to interpret how the undirected features were able to influence the prediction of outcome. Figure 2 visualizes the produced graphical gLIME explanations for the OAI problem including forty features. Specifically, the network of the first permuted dataset (100A) of the 100th observation was depicted. For visualization purposes, very small weight edges ($w < 0.1$) were excluded from the graph in order to identify easier the significant path routes.

The following remarks could be drawn from Figure 2: (i) The classification output (O1) is strongly affected by variables V5 and V33. Green lines reveal a positive effect (an increase in those variables pushes O1 to increase as well) (ii) V5 is also positively affected by V36 that is also similarly correlated with V34. Therefore a path of positive correlations is identified between the variables V34, V36, V5 and the output O1. (iii) Even there is no direct correlation between V34 and the output O1 (as shown in Figure 3 in both LIME and gLIME bar graphs) the produced graph helps the user to identify that there is indirect relationship between them. (iv) There are variables that are strongly correlated to each other but have no (or minor) effect to the output. (v) This descriptive graphical presentation of the feature importance values and their interconnections enhance the users’ understanding of the rationale behind the decision making mechanism of the trained black box model.

Figure 3 shows the barplot explanations of both LIME and gLIME. The most important features that affect the predicted outcome are placed at the top of the list. Blue and orange color show positive and negative influence of each feature with the output feature (O1). Similar ranking are produced by both approaches however due to regularized partial correlations of gLIME, the effect of the last seven features is significantly lowered compared to LIME.

The results of our proposed algorithm (gLIME) are showing superiority in terms of stability and interpretability compared to LIME in tabular data; however, more research on this topic needs to be undertaken. Therefore, we are planning to apply gLIME to other types of data, such as text and images. Moreover, further research should be done to investigate the robustness of the predicted outcomes when more features (over 100) are included. Finally, introducing causality in the produced graphs is within our future scientific interests that will lead to more meaningful directed graphical explanations.

Figure 2: The graphical representation of the first permuted dataset of the 100th observation. For the sake of simplicity, weights under 0.1 were eliminated from the graph.



V1	V00LKPFCRE	V8	V00RKPFCRE	V15	V00SRVDY	V22	V00KSXRKN5	V29	V00KOOSKPL	V36	P01KSX
V2	V00KSXRKN1	V9	V00KQOL1	V16	V00SMKPKYR	V23	V00KQOL3	V30	V00DTLUT	V37	P01KSURGL
V3	V00KSXLKN5	V10	V00ABCIRC	V17	V00RX30NUM	V24	V00KQOL2	V31	V00D1B12	V38	V00rmaxf
V4	V00KSXLKN1	V11	P02KSURG	V18	V00RKPFHDEG	V25	V00KOOSYMR	V32	V00BPDIAS	V39	V00reIHPL
V5	P02ELGRISK	V12	P01OADEGCV	V19	V00PA530	V26	V00KOOSYML	V33	V00AGE	V40	V00IITHRL
V6	P01KSURGR	V13	P01BMI	V20	V00LKRFXPN	V27	V00KOOSQOL	V34	P02KPNLCV	O1	GROUPS
V7	V00WTMAXKG	V14	V00remaxf	V21	V00LKFHDEG	V28	V00KOOSKPR	V35	P01WEIGHT		

4 Conclusions

gLIME is a novel graphical model-agnostic explainability methodology that goes beyond the current state-of-the-art of XAI methods. It incorporates graph theory to identify the most relevant features that affect the black box model’s output, highlight hidden but important relationships among features quantifying the strength of the relationships and determine significant path routes presenting how the information flows from one node to the end node (the predicted output). This enhanced graphical visualization of the produced explanations can increase significantly the user’s understanding of the inner workings and the reasoning behind the ML models’ decision-making process. Apart from being interpretable, gLIME was also proved to be more robust and consistent compared to LIME on an extensive experimentation that included two well-known classification datasets. gLIME’s informative and graphically given explanations that could unlock black boxes contributing to the development of robust and trustworthy AI-empowered systems.

5 Acknowledgement

This research was funded by the European Community’s H2020 Programme, under grant agreement No. 957362 (XMANAI).

References

- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- David Gunning and David Aha. Darpa’s explainable artificial intelligence (xai) program. *AI Magazine*, 40(2):44–58, 2019.
- Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable ai: A review of machine learning interpretability methods. *Entropy*, 23(1):18, 2021.

- Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, pages 3145–3153. PMLR, 2017.
- Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- 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*, pages 1135–1144, 2016.
- Damien Garreau and Ulrike Luxburg. Explaining the explainer: A first theoretical analysis of lime. In *International Conference on Artificial Intelligence and Statistics*, pages 1287–1296. PMLR, 2020.
- Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*, 2017.
- Sacha Epskamp and Eiko I Fried. A tutorial on regularized partial correlation networks. *Psychological methods*, 23(4): 617, 2018.
- Nicolai Meinshausen, Peter Bühlmann, et al. High-dimensional graphs and variable selection with the lasso. *Annals of statistics*, 34(3):1436–1462, 2006.
- Sergei N Dorogovtsev and José FF Mendes. *Evolution of networks: From biological nets to the Internet and WWW*. OUP Oxford, 2013.
- David Hevey. Network analysis: a brief overview and tutorial. *Health Psychology and Behavioral Medicine*, 6(1): 301–328, 2018.
- Steffen L Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Denny Borsboom and Angélique OJ Cramer. Network analysis: an integrative approach to the structure of psychopathology. *Annual review of clinical psychology*, 9:91–121, 2013.
- Richard J. McNally. Can network analysis transform psychopathology? *Behaviour Research and Therapy*, 86:95–104, 2016. ISSN 0005-7967. doi:<https://doi.org/10.1016/j.brat.2016.06.006>. URL <https://www.sciencedirect.com/science/article/pii/S0005796716301103>. Contributions from experimental psychopathology to the understanding and treatment of mental disorders.
- Giulio Costantini, Sacha Epskamp, Denny Borsboom, Marco Perugini, René Möttus, Lourens J Waldorp, and Angélique OJ Cramer. State of the art personality research: A tutorial on network analysis of personality data in r. *Journal of Research in Personality*, 54:13–29, 2015.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Peng Zhao and Bin Yu. On model selection consistency of lasso. *The Journal of Machine Learning Research*, 7: 2541–2563, 2006.
- Mathias Drton and Michael D Perlman. Model selection for gaussian concentration graphs. *Biometrika*, 91(3):591–602, 2004.
- Rina Foygel and Mathias Drton. Extended bayesian information criteria for gaussian graphical models. *arXiv preprint arXiv:1011.6640*, 2010.
- Maurice George Kendall. Rank correlation methods. 1948.

Figure 3: The barplot explanations of the model-agnostic algorithms *gLIME* (left) and LIME (right).

