

# Unwrapping The Black Box of Deep ReLU Networks: Interpretability, Diagnostics, and Simplification

Agus Sudjianto<sup>1</sup>, William Knauth<sup>2</sup>, Rahul Singh<sup>1</sup>  
Zebin Yang<sup>3</sup>, and Aijun Zhang<sup>3,\*</sup>

<sup>1</sup>Wells Fargo, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>The University of Hong Kong

## Abstract

The deep neural networks (DNNs) have achieved great success in learning complex patterns with strong predictive power, but they are often thought of as “black box” models without a sufficient level of transparency and interpretability. It is important to demystify the DNNs with rigorous mathematics and practical tools, especially when they are used for mission-critical applications. This paper aims to unwrap the black box of deep ReLU networks through local linear representation, which utilizes the activation pattern and disentangles the complex network into an equivalent set of local linear models (LLMs). We develop a convenient LLM-based toolkit for interpretability, diagnostics, and simplification of a pre-trained deep ReLU network. We propose the local linear profile plot and other visualization methods for interpretation and diagnostics, and an effective merging strategy for network simplification. The proposed methods are demonstrated by simulation examples, benchmark datasets, and a real case study in home lending credit risk assessment.

**Keywords:** Activation pattern, Deep ReLU network, Local linear interpretability, Network diagnostics, Network simplification.

## 1 Introduction

Explainable artificial intelligence (XAI) or interpretable machine learning (IML) has received great attention in recent years; see [22, 34, 36, 5] among others. The majority of machine learning algorithms have the ability to accurately predict various complex phenomena, but they are often regarded as black-box models since their inner decision-making process cannot

---

\*Corresponding author: ajzhang@umich.edu

be easily understood by human beings. A fundamental objective of XAI/IML is to understand why these black-box models work for different real-world problems. This is important for AI/ML model usage in mission-critical applications, not only due to the increased regulatory scrutiny but also in order to gain the trust of the model users. Without a good understanding of why an AI/ML model works, it may perform abnormally as the data feed is slightly changed, for an instance, the adversarial perturbation in computer vision [20].

There exist two lines of XAI/IML research works about model interpretation, namely intrinsic model interpretability and post-hoc model explainability. Intrinsic interpretability is to study the inherent statistical properties of a specific machine learning model. It is argued by [47] that the intrinsic interpretability of a complex model ought to be induced from practical constraints, including additivity, sparsity, monotonicity, near-orthogonality, linearity, and smoothness. In this regard, some classical models in statistics are actually intrinsically interpretable, including generalized linear model (GLM) [33], generalized additive model (GAM) [24], and generalized additive index model (GAIM, also known as projection pursuit learning) [16, 26]. In the new IML literature, intrinsically interpretable models that integrate deep neural networks include GAM-based NAM [1] and GAMI-Net [48], and GAIM-based xNN [46], ExNN [47] and AxNN [12]. In computer vision, there exist intrinsically interpretable models based on convolutional neural networks [49, 2, 8].

On the other hand, post-hoc explainability is model-agnostic. For any black-box model that produces an output from an input, post-hoc analytical methods include variable importance (VI) [10], partial dependence plot (PDP) [15], individual conditional expectation (ICE) plot [18], and accumulated local effects (ALE) plot [3]. Post-hoc explainability may also rely on surrogate modeling, including both the global surrogates as in model compression or distillation [11, 6, 25] and the perturbation-based local surrogates as in LIME [38] and SHAP [32, 31]. However, these post-hoc methods are often based on crude assumptions, and their interpretation results are not fully reliable. For an instance, [43] recently proposed an adversarial attack method that may hide the model bias from LIME and SHAP.

In this paper, we investigate the intrinsic interpretability of deep neural networks (DNNs) with ReLU (rectified linear unit) activation functions, which are arguably the most popular type of neural networks in deep learning [29, 19]. The deep ReLU networks enjoy many appealing properties, including better performance than sigmoidal activations [17], rich expressivity, and universal approximation ability [35, 30, 4, 39], double-descent risk curve upon over-parameterization [9], fast and scalable training algorithms by stochastic gradient descent (SGD) and its variants [50]. Despite their strong predictive power, the deep ReLU networks act as a black box without a sufficient level of transparency, interpretability, and robustness. The related works on ReLU DNN interpretation are reviewed below.

## 1.1 Related Works

Besides the model-agnostic post-hoc methods reviewed above, there exist some works that tried to interpret pre-trained deep ReLU networks specifically. One type of existing works is to quantify the individual feature importance by computing integrated gradients [44] or decomposing predictions through layer-wise backpropagation [7, 42], as an attempt to attribute the influence on prediction to each input feature in a rank-order manner. There exist similar methods for detecting feature interactions in deep ReLU networks, including pairwise feature interactions by integrated Hessians [27] and general neural interactions by aggregated weights [45]. However, all these methods are based on approximations and they cannot guarantee consistent interpretation results.

An important fact about ReLU DNNs is that the function expressed by a ReLU network is piecewise linear [35]. One popular measure of expressivity is the total number of distinct activation patterns [37, 41, 23], since an activation pattern is known to be associated with a convex polytope of local linear features. The use of activation pattern also appears in [13] (under the name of hidden neuron configuration) in study of piecewise linear neural networks with binary responses, for which a so-called exact interpretability is proposed to interpret the polytope boundary features. Besides, the activation pattern is used recently by [21] as a central tool to explain the neural network behavior in terms of input and layer properties.

For deep ReLU networks that are usually over-parametrized, pruning is an effective way to simplify the network; see e.g. [14] with “lottery ticket hypothesis” and [40] based on mixed-integer linear programming. For a large network, pruning is to fix certain entries of weight matrices to be zero, resulting in subnetworks with fewer unknown parameters. It is evident from multiple recent works that such pruned subnetworks may not only predict well without sacrifice on prediction performance but also enjoy the improved interpretability due to network simplification. The interpretability or property verification of deep ReLU networks based on network simplification can be also referred to as [21, 28].

## 1.2 Main Contributions

Our purpose of this study is to unwrap the black box of ReLU DNNs and endue them with intrinsic interpretability. We develop a convenient toolkit called **Aletheia**<sup>©</sup> for interpretability, diagnostics, and simplification of deep ReLU networks. Our main contributions in this paper are summarized as follows.

1. An unwrapper is developed to disentangle a pre-trained ReLU DNN into an equivalent set of local linear models (LLMs). Such unwrapper for both regression and classification problems is rigorously supported by local linear representation of ReLU DNNs based on the device of activation pattern and activation region.

2. LLM-based intrinsic interpretability is developed for deep ReLU networks, including local linear profile and joint importance of an individual feature, simultaneous interpretation of multiple LLMs through the parallel coordinate plot, and local interpretability by statistical inference.
3. LLM-based diagnostics are developed for deep ReLU networks for checking the sizes of local regions, the extrapolation strengths of LLMs, and the similarities among the LLM coefficients. Single-instance or single-class regions are found to be a critical issue for over parametrized ReLU DNNs.
4. An effective merging algorithm is proposed for LLM-based network simplification. It combines the locally homogeneous LLMs with similar coefficients subject to local region connectivity. The single-instance or single-class regions are merged to nearest neighbor LLMs with large/medium sample sizes. Besides, an optional flattening strategy is proposed to further simplify the network.
5. A new Python package called **Aletheia**<sup>©</sup> is developed as a convenient toolkit for unwrapping the black box of deep ReLU networks. It includes efficient algorithms for local linear representation and network simplification, as well as visualization tools for interpretation and diagnostics. The use of **Aletheia**<sup>©</sup> is demonstrated through simulation examples, benchmark datasets, and a real case study in credit risk assessment.

The rest of the paper is organized as follows. In Section 2 we derive the local linear representation for deep ReLU networks based on the activation pattern and activation region. Then Section 3 presents the LLM-based interpretability, and Section 4 discusses LLM-based diagnostics. In Section 5 we develop the merging and flattening techniques for network simplification. Section 6 presents a real case study in home lending credit risk assessment, as an illustration of the proposed methodology in real-life environments. Section 7 wraps up the paper with concluding remarks.

## 2 Local Linear Representation

Denote by  $\mathcal{N}$  a feedforward ReLU network that has the input layer with  $\mathbf{x} \in \Omega \subseteq \mathbb{R}^d$  and  $L$  hidden layers with neuron sizes  $[n_1, \dots, n_L]$ . For each hidden neuron  $\mathbf{u}_i^{(l)}$ , denote its input (from left) as  $z_i^{(l)}$  and the output (to right) as  $\chi_i^{(l)}$ , then by the ReLU activation

$$\chi_i^{(l)} = \max\{0, z_i^{(l)}\}, \quad \text{for } i = 1, \dots, n_l \text{ and } l = 1, \dots, L. \quad (1)$$

From layer  $l-1$  to  $l$ , given the weight matrix  $\mathbf{W}^{(l-1)}$  of size  $n_l \times n_{l-1}$  and the bias vector  $\mathbf{b}^{(l-1)}$  of length  $n_l$ , the layer- $(l-1)$  output  $\boldsymbol{\chi}^{(l-1)}$  leads to the layer- $l$  input  $\mathbf{z}^{(l)}$  by the affine

transformation

$$\mathbf{z}^{(l)} = \mathbf{W}^{(l-1)}\boldsymbol{\chi}^{(l-1)} + \mathbf{b}^{(l-1)}, \quad \text{for } l = 1, \dots, L, \quad (2)$$

where the layer zero corresponds to the input layer with  $\boldsymbol{\chi}^{(0)} = \mathbf{x}$ ,  $\mathbf{W}^{(0)}$  of size  $n_1 \times d$ , and  $\mathbf{b}^{(0)}$  of length  $n_1$ . As for the output layer (i.e. layer  $L + 1$ ), the layer- $L$  output  $\chi_i^{(L)}$  can be used to predict the final response  $y$  by the generalized linear model (GLM),

$$\mathbb{E}[y] = \sigma(\mathbf{w}^{(L)}\boldsymbol{\chi}^{(L)} + b^{(L)}) \equiv \sigma(\eta(\mathbf{x})) \quad (3)$$

where the  $\sigma$  function can be the identity link for a regression problem (i.e.  $y$  is continuous) or the logit link for a classification problem (i.e.  $y$  is binary). In such GLM formulation with univariate response,  $\mathbf{W}^{(L)} \equiv \mathbf{w}^{(L)}$  is a row vector and  $\mathbf{b}^{(L)} \equiv b^{(L)}$  is a scalar. For convenience, let us use  $\boldsymbol{\theta}$  to denote the set of network parameters  $\{\mathbf{W}^{(l)}, \mathbf{b}^{(l)}, l = 0, 1, \dots, L\}$ . See Figure 1 for a toy ReLU net with bivariate inputs, two hidden layers with node sizes  $[2, 4]$ , and a univariate output. Note that when  $y$  is one-hot encoded (either binary or multi-category classification), we can use the softmax link with a straightforward extension.

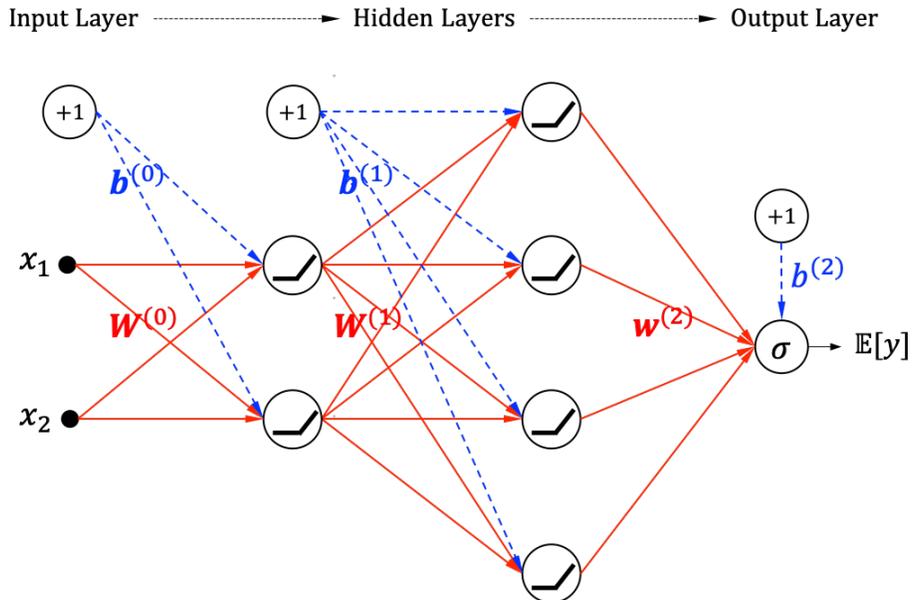


Figure 1: Toy ReLU Net with two hidden layers of node sizes  $[2, 4]$ . The network parameters  $\boldsymbol{\theta}$  include the weights  $\mathbf{W}_{2 \times 2}^{(0)}$ ,  $\mathbf{W}_{4 \times 2}^{(1)}$ ,  $\mathbf{w}_{1 \times 4}^{(2)}$  and the biases  $b_{2 \times 1}^{(0)}$ ,  $b_{4 \times 1}^{(1)}$ ,  $b_{1 \times 1}^{(2)}$ .

## 2.1 Activation Pattern

For each hidden neuron, the ReLU activation (1) has two states, namely “on” (active, when  $z \geq 0$ ) and “off” (inactive, when  $z < 0$ ). To investigate the on/off states of all hidden neurons in the ReLU network, we employ the notion of activation pattern [37] and formally define it as follows.

**Definition 1 (Activation Pattern)** For  $\mathcal{N}$  with  $L$  hidden layer sizes  $[n_1, \dots, n_L]$ , define the activation pattern to be the following binary vector

$$\mathbf{P} = [\mathbf{P}^{(1)}; \dots; \mathbf{P}^{(L)}] \in \{0, 1\}^{\sum_{i=1}^L n_i} \quad (4)$$

which indicates the on/off state of each hidden neuron in the network. Specifically, the component  $\mathbf{P}^{(l)}$  is called a layered pattern for  $l = 1, \dots, L$ . The activation pattern  $\mathbf{P}$  is said to be trivial if there is at least one  $\mathbf{P}^{(l)} \equiv 0$  for some  $l$ .

The length of activation pattern is  $\sum_{i=1}^L n_i$ , which equals the total number of hidden neurons in the network. For each input  $\mathbf{x} \in \Omega$ , through the feedforward neural network with fixed parameters, it would determine the values of  $\{\mathbf{z}^{(l)}, \boldsymbol{\chi}^{(l)}\}$  by (2) and (1) for  $l = 1, \dots, L$  sequentially. Each  $\mathbf{z}^{(l)}$  would determine the on/off states of the ReLU neurons at layer  $l$ , therefore determine a layered pattern, as denoted by  $\mathbf{P}^{(l)}(\mathbf{x})$ . Thus, any input instance  $\mathbf{x}$  corresponds to a particular activation pattern of the form

$$\mathbf{P}(\mathbf{x}) = [\mathbf{P}^{(1)}(\mathbf{x}); \dots; \mathbf{P}^{(L)}(\mathbf{x})], \quad (5)$$

which can be either trivial or non-trivial. We are now ready to define the *activation region* in association with each distinct pattern

$$\mathcal{R}^{\mathbf{P}} = \{\mathbf{x} \in \Omega : \mathbf{P}(\mathbf{x}) = \mathbf{P}\}, \quad (6)$$

which is known to be a convex polytope with closed-form boundaries to be discussed below.

The activation region plays an important role in understanding how the hidden layers/nodes partition the full space  $\Omega$  into disjoint sub-spaces. To illustrate how it works, consider the toy ReLU net in Figure 1 with the following tailor-made weight and bias parameters for the two hidden layers,

$$\mathbf{W}^{(0)} = \frac{1}{\sqrt{2}} \begin{pmatrix} -1 & 1 \\ 1 & 1 \end{pmatrix}, \quad \mathbf{b}^{(0)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}; \quad \mathbf{W}^{(1)} = \begin{pmatrix} 1 & 1/4 \\ 1/2 & 1/3 \\ 1/3 & 1/2 \\ 1/4 & 1 \end{pmatrix}, \quad \mathbf{b}^{(1)} = \frac{3}{10} \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix}. \quad (7)$$

In this example, the affine transformation  $\mathbf{z}^{(1)} = \mathbf{W}^{(0)}\mathbf{x} + \mathbf{b}^{(0)}$  is an orthogonal rotation. It is easy to check the boundary conditions between on/off states of two hidden nodes in the first layer, given by

$$\begin{cases} \text{on: } -x_1 + x_2 > 0 \\ \text{off: } -x_1 + x_2 \leq 0 \end{cases} \quad \text{and} \quad \begin{cases} \text{on: } x_1 + x_2 > 0 \\ \text{off: } x_1 + x_2 \leq 0 \end{cases}.$$

See Figure 2 (left) for the plot of such boundaries within the 2D domain  $[-1, 1]^2$ . The domain is divided into four sub-domains, each corresponding to a unique activation pattern

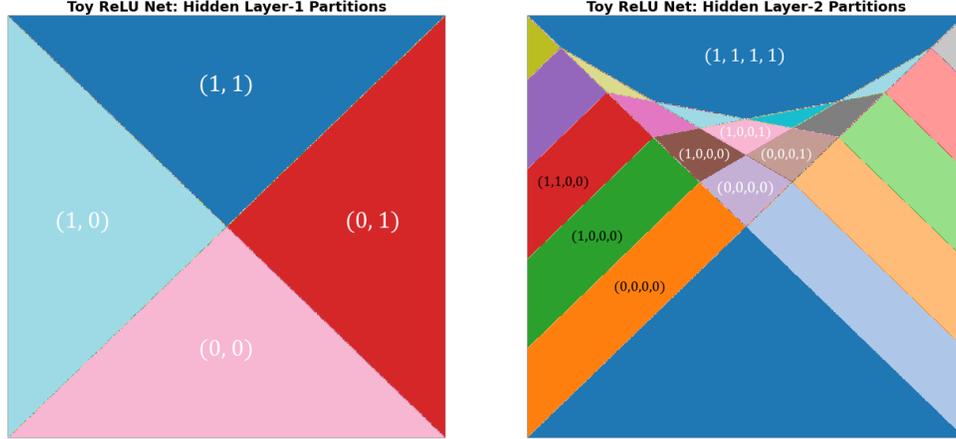


Figure 2: Region partitions by the toy ReLU net with tailor-made parameters (7).

(as annotated). More specifically, it includes three non-trivial patterns  $(1,1)$ ,  $(1,0)$ ,  $(0,1)$  and one trivial pattern  $(0,0)$ .

Now consider the forward propagation from the first hidden layer to the second hidden layer. Among the four regions shown in Figure 2 (left), the bottom trivial  $(0,0)$ -quadrant receives constant inputs  $\mathbf{z}^{(2)}$ , by  $\boldsymbol{\chi}^{(1)} = \mathbf{0}$  and equation (2); so this region generates the same constant output thereafter. On the other hand, for the top  $(1,1)$ -quadrant, the input is updated to be  $\boldsymbol{\chi}^{(1)} = \mathbf{z}^{(1)}$ . Within this region, the internal boundary conditions between on/off states of four hidden nodes can be determined by the affine-transformed variables  $\mathbf{z}^{(2)} = \mathbf{W}^{(1)}\mathbf{z}^{(1)} + \mathbf{b}^{(1)}$ , i.e.

$$\begin{aligned}
 \text{Node 1: } & z_1^{(1)} + \frac{1}{4}z_2^{(1)} - \frac{3}{10} > 0 \text{ or } \leq 0 \\
 \text{Node 2: } & \frac{1}{2}z_1^{(1)} + \frac{1}{3}z_2^{(1)} - \frac{3}{10} > 0 \text{ or } \leq 0 \\
 \text{Node 3: } & \frac{1}{3}z_1^{(1)} + \frac{1}{2}z_2^{(1)} - \frac{3}{10} > 0 \text{ or } \leq 0 \\
 \text{Node 4: } & \frac{1}{4}z_1^{(1)} + z_2^{(1)} - \frac{3}{10} > 0 \text{ or } \leq 0
 \end{aligned}$$

These four inequalities divide the domain into 11 sub-regions in total, as visualized in the top quadrant of Figure 2 (right)<sup>1</sup>. For each of the 11 sub-regions, its corresponding activation pattern can be derived from the on/off states of hidden nodes. Due to small space sizes, only five sub-regions are annotated (in white-colored text) with explicit activation patterns. Similarly, we can analyze the left  $(1,0)$ -quadrant and the right  $(0,1)$ -quadrant and obtain

<sup>1</sup>Note that this visualization is similar to Figure 2 in [41] who used a ReLU net with 3 hidden layers of  $[2,2,2]$  nodes. However, there is subtle difference in the piecewise linear cuts. Both using 6 hidden nodes, [41] obtained 20 regions in total (which is less than the full potential expressivity for a depth-4 network), and we obtained 22 regions in total (which achieves the full potential expressivity for a depth-3 network).

their partitions by the four hidden nodes. In the above inequalities, muting  $z_2^{(1)}$  for the  $(1,0)$ -region (or muting  $z_1^{(1)}$  for the  $(0,1)$ -region) leads to the internal boundary conditions, as well as the corresponding activation patterns. In Figure 2 (right), the left  $(1,0)$ -quadrant is divided into 5 sub-regions, for which we annotate (in black-colored text) three examples of layered activation patterns.

By this toy example, it is clear that the original input space is sequentially partitioned into multiple activation regions, and each region is associated with a unique activation pattern of the hidden nodes in each layer. It is also obvious that each activation region is a convex polytope, since all the boundary conditions are linear inequalities. In the meanwhile, it is important to observe from this example that the toy ReLU net generates the activation patterns with certain rules; in another word, an arbitrary binary pattern is not always expressible. For an instance,  $(1,0;0,0,1,1)$  is not a legitimate activation pattern in Figure 2.

For ReLU DNNs in general, we have the following lemma that has partially appeared before in [37, 23]. We use  $\mathcal{P}_{\text{expr}}(\mathcal{N})$  to denote the set of distinct activation patterns that are expressible by a ReLU net  $\mathcal{N}$ .

**Lemma 1 (Convex Partitioning)** *For a ReLU DNN  $\mathcal{N}$ , each activation pattern  $\mathbf{P} = [\mathbf{P}^{(1)}; \dots; \mathbf{P}^{(L)}]$  in  $\mathcal{P}_{\text{expr}}(\mathcal{N})$  is associated with a convex activation region  $\mathcal{R}^{\mathbf{P}} \subset \Omega$ , satisfying the following inequality constraints based on the layer-wise affine-transformed variables:*

$$(-1)^{\mathbf{P}^{(l)}} \circ \mathbf{z}^{(l)} \leq \mathbf{0}, \quad l = 1, \dots, L \quad (8)$$

where  $\circ$  denotes the Hadamard product (element-wise). Moreover, for any two distinct activation patterns, their corresponding activation regions are disjoint. Hence, the ReLU net  $\mathcal{N}$  partitions the original input space into a finite set of convex sub-spaces, i.e.,

$$\Omega = \bigcup_{\mathbf{P} \in \mathcal{P}_{\text{expr}}(\mathcal{N})} \mathcal{R}^{\mathbf{P}}. \quad (9)$$

In Lemma 1,  $\mathcal{P}_{\text{expr}}(\mathcal{N})$  is yet to be determined. This has been an active topic in the expressivity study of DNNs, in particular the upper and lower bounds of the total number of distinct patterns [35, 37, 41, 23]. According to [37], the cardinality of  $\mathcal{P}_{\text{expr}}(\mathcal{N})$  can be as large as  $O(k^{dL})$  for a ReLU net with  $L$  hidden layers of width  $k$ . In this paper, our interest is to interpret a pre-trained network  $\mathcal{N}_{\text{train}}$  based on training dataset  $\mathcal{X}_{\text{train}}$ , for which the useful activation patterns (to be unwrapped in Section 2.3) turns out to be much fewer.

## 2.2 Local Linear Models

A ReLU DNN partitions the input space into many sub-regions, and each sub-region is determined by a unique activation pattern. By multi-layer propagation, the original input

features are sequentially transformed by (2) subject to the ReLU activation (1). Finally, the prediction takes the linear form

$$\eta(\mathbf{x}) = \mathbf{w}^{(L)} \boldsymbol{\chi}^{(L)} + \mathbf{b}^{(L)}, \quad (10)$$

up to a pre-specified  $\sigma$  link function. In this section, we derive the final features  $\boldsymbol{\chi}^{(L)}$  explicitly for each sub-region and obtain the local linear models in closed form.

For a trivial activation pattern when at least one layer has all inactive neurons, the output of such a layer would become all zero. For the corresponding trivial regions (e.g. the sub-regions annotated with (0,0) and (0,0,0,0)'s in Figure 2), the prediction would be a constant value within each region, since the final output is only affected by the layer-wise bias terms, but not the original input variables  $\mathbf{x}$ .

For a non-trivial activation pattern  $\mathbf{P} = [\mathbf{P}^{(1)}; \dots; \mathbf{P}^{(L)}]$ , let us introduce a diagonal matrix  $\mathbf{D}^{(l)}$  for each layer, whose diagonal takes the same (0,1)-entries as  $\mathbf{P}^{(l)}$ , i.e.,

$$\mathbf{D}^{(l)} = \text{diag}(\mathbf{P}^{(l)}), \quad \text{for } l = 1, \dots, L.$$

Using this notation, the layer-wise output after ReLU activation (1) can be effectively expressed in vector form as the following chain rule,

$$\boldsymbol{\chi}^{(l)} = \max\{\mathbf{0}, \mathbf{z}^{(l)}\} = \mathbf{D}^{(l)} \mathbf{z}^{(l)} = \mathbf{D}^{(l)} (\mathbf{W}^{(l-1)} \boldsymbol{\chi}^{(l-1)} + \mathbf{b}^{(l-1)}). \quad (11)$$

Thus, we may derive the final features  $\boldsymbol{\chi}^{(L)}$  recursively as follows,

$$\begin{aligned} \boldsymbol{\chi}^{(L)} &= \mathbf{D}^{(L)} (\mathbf{W}^{(L-1)} \boldsymbol{\chi}^{(L-1)} + \mathbf{b}^{(L-1)}) \\ &= \mathbf{D}^{(L)} \mathbf{W}^{(L-1)} \mathbf{D}^{(L-1)} (\mathbf{W}^{(L-2)} \boldsymbol{\chi}^{(L-2)} + \mathbf{b}^{(L-2)}) + \mathbf{D}^{(L)} \mathbf{b}^{(L-1)} \\ &= \dots \\ &= \mathbf{D}^{(L)} \mathbf{W}^{(L-1)} \dots \mathbf{D}^{(1)} \mathbf{W}^{(0)} \mathbf{x} + \sum_{l=1}^{L-1} \mathbf{D}^{(L)} \mathbf{W}^{(L-1)} \dots \mathbf{D}^{(l+1)} \mathbf{W}^{(l)} \mathbf{D}^{(l)} \mathbf{b}^{(l-1)} + \mathbf{D}^{(L)} \mathbf{b}^{(L-1)} \\ &= \prod_{h=1}^L \mathbf{D}^{(L+1-h)} \mathbf{W}^{(L-h)} \mathbf{x} + \sum_{l=1}^{L-1} \prod_{h=1}^{L-l} \mathbf{D}^{(L+1-h)} \mathbf{W}^{(L-h)} \mathbf{D}^{(l)} \mathbf{b}^{(l-1)} + \mathbf{D}^{(L)} \mathbf{b}^{(L-1)}. \end{aligned}$$

Furthermore, the linear prediction (10) by the output layer is given explicitly by

$$\eta(\mathbf{x}) = \prod_{h=1}^L \mathbf{W}^{(L+1-h)} \mathbf{D}^{(L+1-h)} \mathbf{W}^{(0)} \mathbf{x} + \sum_{l=1}^L \prod_{h=1}^{L+1-l} \mathbf{W}^{(L+1-h)} \mathbf{D}^{(L+1-h)} \mathbf{b}^{(l-1)} + \mathbf{b}^{(L)},$$

in which  $\mathbf{W}^{(L)} \equiv \mathbf{w}^{(L)}$  for the notation convenience. Hence, we have the following theorem that characterizes deep ReLU networks via local linear models (LLMs).

**Theorem 1 (Local Linear Model)** *For a ReLU DNN and any of its expressible activation pattern  $\mathbf{P}$ , the local linear model on the activation region  $\mathcal{R}^{\mathbf{P}}$  is given by*

$$\eta^{\mathbf{P}}(\mathbf{x}) = \tilde{\mathbf{w}}^{\mathbf{P}} \mathbf{x} + \tilde{\mathbf{b}}^{\mathbf{P}}, \quad \forall \mathbf{x} \in \mathcal{R}^{\mathbf{P}} \quad (12)$$

with the following closed-form parameters

$$\tilde{\mathbf{w}}^{\mathbf{P}} = \prod_{h=1}^L \mathbf{W}^{(L+1-h)} \mathbf{D}^{(L+1-h)} \mathbf{W}^{(0)}, \quad \tilde{\mathbf{b}}^{\mathbf{P}} = \sum_{l=1}^L \prod_{h=1}^{L+1-l} \mathbf{W}^{(L+1-h)} \mathbf{D}^{(L+1-h)} \mathbf{b}^{(l-1)} + \mathbf{b}^{(L)}. \quad (13)$$

By combining the region partitioning in Lemma 1 and the local linear models in Theorem 1, we have a local linear representation for ReLU DNNs. In another word, a deep ReLU network can be equivalently re-expressed as a finite set of local linear models each functioning exclusively on a disjoint convex sub-region. It is clear that the activation pattern plays a critical role in such local linear representation. A remaining question is how to determine the activation patterns for network interpretation purpose, for which we develop an effective unwrapper for pre-trained ReLU DNNs.

## 2.3 Unwrapping Pre-trained DNNs

For a generic ReLU net  $\mathcal{N}$ , the number of expressible activation patterns as in  $\mathcal{P}_{\text{expr}}(\mathcal{N})$  grows exponentially in network depth, width and input dimensionality. It is just unrealistic to investigate each possible activation pattern, and unnecessary, too. In this paper, we concentrate on the pre-trained DNNs  $\mathcal{N}_{\text{train}}$  based on training data  $\mathcal{X}_{\text{train}}$ , for which we define the set of training data induced activation patterns:

$$\mathcal{P}_{\text{train}} = \{\mathbf{P} \in \mathcal{P}_{\text{expr}} : \mathcal{R}^{\mathbf{P}} \cap \mathcal{X}_{\text{train}} \neq \emptyset\}, \quad (14)$$

and similarly we can define the set of testing data induced activation patterns  $\mathcal{P}_{\text{test}} = \{\mathbf{P} \in \mathcal{P}_{\text{expr}} : \mathcal{R}^{\mathbf{P}} \cap \mathcal{X}_{\text{test}} \neq \emptyset\}$ . It simply means there is at least one training (resp. testing) instance that will be covered by the  $\mathbf{P}$ -associated activation region for  $\mathbf{P} \in \mathcal{P}_{\text{train}}$  (resp.  $\mathbf{P} \in \mathcal{P}_{\text{test}}$ ). Their relationship is depicted in Figure 3. Meanwhile, denote the *activation set* of training instances for each activation pattern by

$$\mathcal{R}_{\text{train}}^{\mathbf{P}} = \{\mathbf{x} \in \mathcal{X}_{\text{train}} : \mathbf{P}(\mathbf{x}) = \mathbf{P}\}, \quad \text{for } \mathbf{P} \in \mathcal{P}_{\text{train}} \quad (15)$$

and similarly the activation set of testing instances by  $\mathcal{R}_{\text{test}}^{\mathbf{P}}$ . Obviously,  $\mathcal{R}_{\text{train}}^{\mathbf{P}} = \emptyset$  if  $\mathbf{P} \notin \mathcal{P}_{\text{train}}$ . Note that although  $\mathcal{R}_{\text{train}}^{\mathbf{P}}$  can be viewed as a subset of the convex region  $\mathcal{R}^{\mathbf{P}}$ , the set  $\mathcal{R}_{\text{train}}^{\mathbf{P}}$  is a discrete set of instances without formal convex property.

By Theorem 1, the pre-trained ReLU net  $\mathcal{N}_{\text{train}}$  for the training instances can be equivalently represented by the following set of LLMs,

$$\mathcal{L}_{\text{train}} = \{\eta^{\mathbf{P}}(\mathbf{x}) = \tilde{\mathbf{w}}^{\mathbf{P}} \mathbf{x} + \tilde{\mathbf{b}}^{\mathbf{P}}, \quad \text{for } \mathbf{x} \in \mathcal{R}_{\text{train}}^{\mathbf{P}}, \mathbf{P} \in \mathcal{P}_{\text{train}}\} \quad (16)$$

where the LLM coefficients  $\{\tilde{\mathbf{w}}^{\mathbf{P}}, \tilde{\mathbf{b}}^{\mathbf{P}}\}$  are given by (13). Such a formal procedure for determining  $\{\mathcal{L}_{\text{train}}, \mathcal{P}_{\text{train}}, \{\mathcal{R}_{\text{train}}^{\mathbf{P}} : \mathbf{P} \in \mathcal{P}_{\text{train}}\}\}$  is called an **unwrapper** for pre-trained ReLU

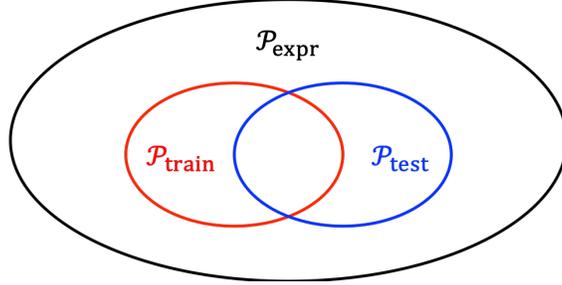


Figure 3: Sets of activation patterns: theoretically expressible  $\mathcal{P}_{\text{expr}}$ , training data induced  $\mathcal{P}_{\text{train}}$  and testing data induced  $\mathcal{P}_{\text{test}}$ .

DNNs. Algorithmically, the unwrapper may identify distinct activation patterns in  $\mathcal{P}_{\text{train}}$  by individual training instances<sup>2</sup>, for which all  $\sum_{l=1}^L n_l$  hidden nodes need to be checked by layer-wise propagated  $\mathbf{z}^{(l)}$  for each  $\mathbf{x} \in \mathcal{X}_{\text{train}}$ . In this way, each identified activation pattern  $\mathbf{P}$  comes with its supporting instances, i.e.  $\mathcal{R}_{\text{train}}^{\mathbf{P}}$ . Finally, each LLM in  $\mathcal{L}_{\text{train}}$  can be evaluated for each unique  $\mathbf{P} \in \mathcal{P}_{\text{train}}$  straightforwardly.

A new Python package called **Aletheia**<sup>©</sup> is under our development, which includes **UnwrapperClassifier** and **UnwrapperRegressor** of ReLU nets for both classification and regression problems. This unwrapper provides a solid foundation for local linear interpretability, diagnostics, and simplification of deep ReLU networks in the next sections.

### 3 LLM-based Network Interpretability

The local linear representation by Theorem 1 enables us to directly interpret ReLU DNNs based on LLMs. In this section, we first introduce the local linear profile for the purpose of interpreting individual feature importances and partial relationships with final prediction, then discuss how to interpret the joint effects of LLMs through effective use of parallel coordinate plot.

For better illustration of the proposed methods, we use the following one-, two- and multi-dimensional examples under both regression and classification settings:

- **ChirpWave** (regression):  $y_i = \sin(2\pi/(x_i + 0.2)) + \varepsilon$  for  $x_i \in [0, 1]$ ,  $i = 1, \dots, 2000$ , where the i.i.d. noise is simulated by  $\varepsilon \sim N(0, 0.1^2)$ ;
- **CoCircles** (binary classification): concentric circle dataset with two features, which can be generated by `sklearn.datasets.make_circles(n_samples=2000, noise=0.1)`;

---

<sup>2</sup>For large-scale training data (millions or more sample sizes), certain smart sampling strategy can be employed for computational speed up. This is currently under our investigation.

- **BostonHouse** (regression): Boston house prices dataset with 506 samples and 13 features, which can be loaded by `sklearn.datasets.load_load_boston()`;
- **FicoHeloc** (binary classification): 9871 samples with 35 features after preprocessing the raw data from the FICO website<sup>3</sup>.

Each dataset is split into training (80%) and testing (20%) sets upon random permutation. For simplicity, we use `sklearn.neural_network` (either `MLPRegressor` or `MLPClassifier`) for implementing the ReLU net with four hidden layers each with 40 neurons, except for the small BostonHouse data with two hidden layers each with 20 neurons. Each ReLU net is trained for 2000 epochs using the Adam optimizer, together with an early stopping rule when the validation performance (20% of training data) does not get improved for 100 epochs.

Table 1: Four datasets and deep ReLU network model performances. Each dataset is split into training (80%) and testing (20%) sets. The ReLU net size stands for the hidden layer architecture. The number of LLMs stands for the cardinality of unwrapped  $\mathcal{C}_{\text{train}}$  from the fitted  $\mathcal{N}_{\text{train}}$  based on the training data. The training and testing performances stand for MSE (for ChirpWave and BostonHouse) and AUC (CoCircles and FicoHeloc).

Dataset	Data Size	ReLU Net Size	Num. LLMs	Training	Testing
ChirpWave	2000 x 1	[40]*4	64	0.0113	0.0117
CoCircles	2000 x 2	[40]*4	530	0.9215	0.9320
BostonHouse	506 x 13	[20]*2	155	0.0022	0.0094
FicoHeloc	9871 x 35	[40]*4	7423	0.8325	0.7980

By the unwrapper introduced in the previous section, each trained ReLU net  $\mathcal{N}_{\text{train}}$  can be re-expressed as a set of LLMs  $\mathcal{L}_{\text{train}}$  together with the corresponding set of activation patterns  $\mathcal{P}_{\text{train}}$ . Note that the total number of LLMs (i.e. cardinality of  $\mathcal{P}_{\text{train}}$ ) depends not only on the training data but also on the local optimal solution trained by SGD optimization. In the meanwhile, we obtain the training and testing performances (MSE for regression problems and AUC for classification problems). These details are listed in Table 1.

### 3.1 Local Linear Profile and Joint Importance

To study the marginal effect of each dimensional feature, we define the local linear profile and the joint importance measure based on the unwrapped LLMs.

<sup>3</sup>FICO 2018 xML challenge: <https://community.fico.com/s/explainable-machine-learning-challenge>

**Definition 2 (Local Linear Profile)** For a pre-trained  $\mathcal{N}_{\text{train}}$  re-expressed by the set of LLMs (16), we define the local linear profile of each dimensional feature  $x_j$  ( $j = 1, \dots, d$ ) by the set of local marginal linear functions upon centering

$$\left\{ \eta^{\mathbf{P}}(x) = \tilde{w}_j^{\mathbf{P}} x - \frac{1}{|\mathcal{R}_{\text{train}}^{\mathbf{P}}|} \sum_{\mathbf{x}_i \in \mathcal{R}_{\text{train}}^{\mathbf{P}}} \tilde{w}_j^{\mathbf{P}} x_{ij}, \mathbf{P} \in \mathcal{P}_{\text{train}} \right\}, \quad (17)$$

and define the local intercept profile by the set of local bias terms  $\{\tilde{b}^{\mathbf{P}}, \mathbf{P} \in \mathcal{P}_{\text{train}}\}$ . Furthermore, we define the joint importance of intercept and individual features by

$$\begin{aligned} \text{JI}(b_0) &= \frac{\sum_{\mathbf{P} \in \mathcal{P}_{\text{train}}} |\mathcal{R}_{\text{train}}^{\mathbf{P}}| (\tilde{b}^{\mathbf{P}})^2}{T} \quad (b_0: \text{intercept}) \\ \text{JI}(x_j) &= \frac{\sum_{\mathbf{P} \in \mathcal{P}_{\text{train}}} |\mathcal{R}_{\text{train}}^{\mathbf{P}}| (\tilde{w}_j^{\mathbf{P}})^2}{T}, \quad j = 1, \dots, d \end{aligned} \quad (18)$$

in which  $T = \sum_{\mathbf{P} \in \mathcal{P}_{\text{train}}} |\mathcal{R}_{\text{train}}^{\mathbf{P}}| (\tilde{b}^{\mathbf{P}})^2 + \sum_{j=1}^d \sum_{\mathbf{P} \in \mathcal{P}_{\text{train}}} |\mathcal{R}_{\text{train}}^{\mathbf{P}}| (\tilde{w}_j^{\mathbf{P}})^2$ .

By Definition 2, it is straightforward to check the partial dependence of DNN prediction on each individual feature. We introduce a local linear profile plot, which visualizes the top marginal LLMs (ranked by the cardinality  $|\mathcal{R}_{\text{train}}^{\mathbf{P}}|$ ) from (17), together with the sample distributions of the corresponding regions. Such a profile plot consists of multiple line segments each representing a marginal LLM upon centering. For the ChirpWave example, the DNN-predicted curve by the top-30 LLMs is shown in Figure 4 (left), and the corresponding profile plot of their marginal LLMs is in Figure 4 (right). Each line segment in the profile plot is supported by a small region without overlapping since all the marginal regions are disjoint in this one-dimensional case. Note that in the bottom part of the profile plot, the sample distribution for each activation region is separately smoothed by kernel density estimation upon normalization; so it does not reflect the total number of instances per region.

For the two-dimensional CoCircles example, the LLMs can also be visualized on a plane. See Figure 5 (left) for the resulting 500+ LLMs unwrapped from the ReLU net. It is clear that these LLMs altogether approximate the underlying circle decision boundary in a piecewise linear manner. Each LLM corresponds to an activation region in the input space; as shown in Figure 5 (middle). In this visualization, the large-connected blue regions represents the uncharted domain without training instances, for which no LLMs are unwrapped. The local linear profile plot (with top 30 marginal LLMs) for one of the features is shown in Figure 5 (right). In this case, the projections of local activation regions to a marginal feature may overlap with each other, and the marginal LLMs may show the cross-over pattern.

For datasets with multi-dimensional features, according to Definition 2, we can quantify the joint importance of each feature and rank them in the decreasing order. See the bar charts of the sorted joint importance values for the BostonHouse and FicoHeloc examples

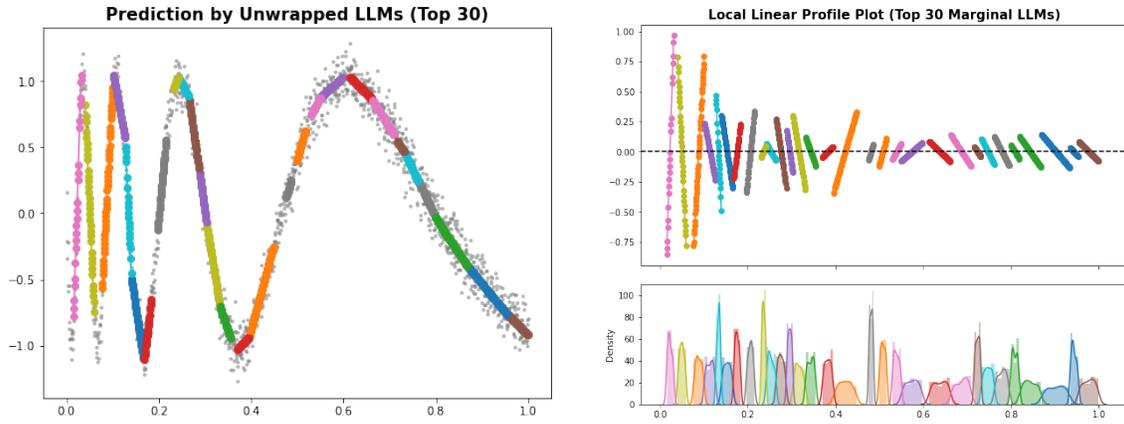


Figure 4: DNN prediction and local linear profile plot (Data: ChirpWave; ReLU Net:  $[40]^*4$ )

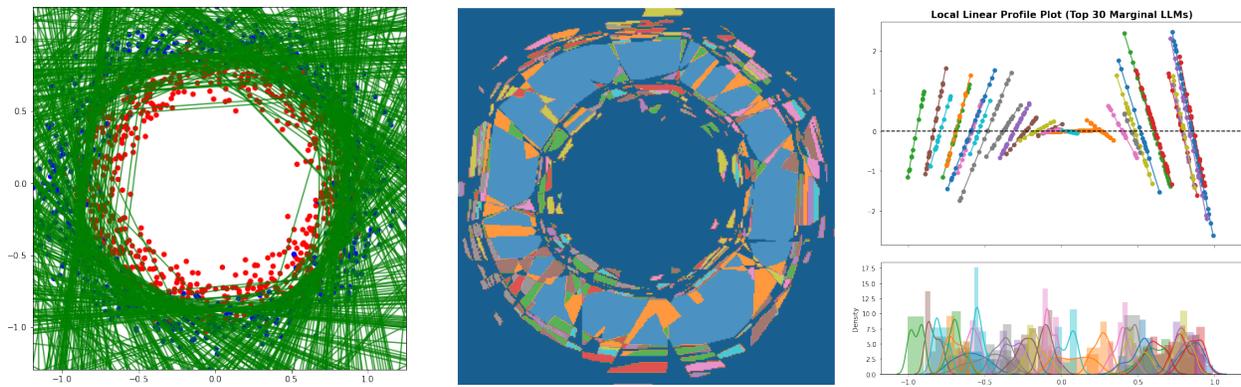


Figure 5: Local linear models unwrapped from pre-trained ReLU net, the corresponding activation regions, and the local linear profile plot (Data: CoCircles; ReLU Net:  $[40]^*4$ )

in Figures 6–7 (left). After that, we can use the local linear profile plot to check the partial dependence for the identified important features. Figures 6–7 (right) show the profile plots of the most important feature based on the top-10 LLMs in each example.

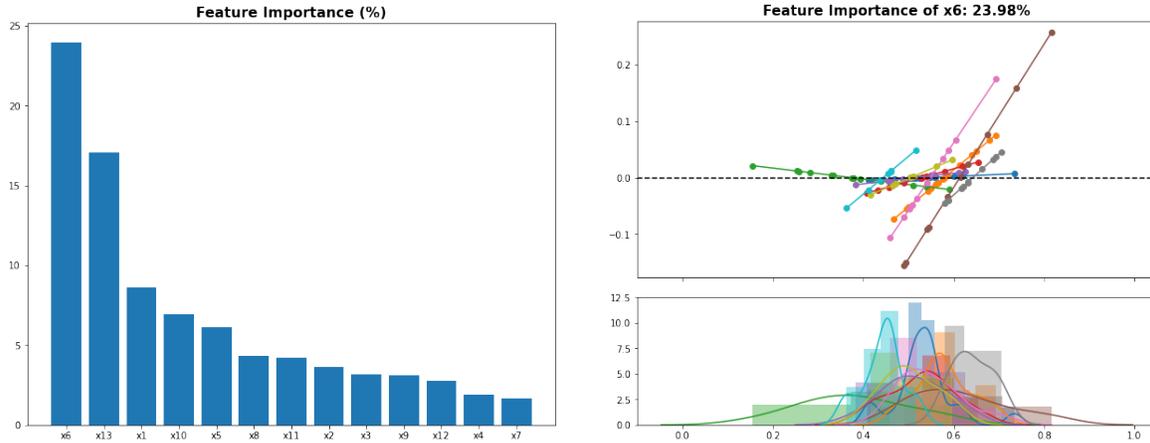


Figure 6: Left: feature importance plot. Right: profile plot of the most important feature. (Data: BostonHouse; ReLU Net: [20]\*2)

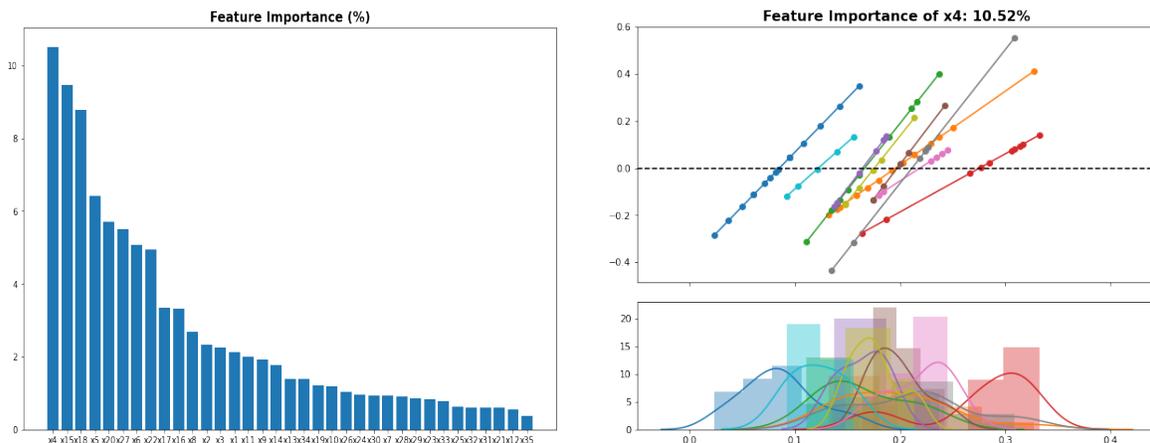


Figure 7: Left: feature importance plot. Right: profile plot of the most important feature. (Data: FicoHeloc; ReLU Net: [40]\*4)

Thus, we have developed both the feature importance and profile plot for the main effect interpretation of pre-trained deep ReLU networks. This is similar to the post-hoc VI and PDP methods developed by [10, 15]. However, as a key difference, our proposed method for ReLU DNNs is an intrinsic interpretability approach based on the unwrapped LLMs.

## 3.2 Parallel Coordinate Plot

It is also important to interpret the joint effects of multiple features simultaneously, which can be either difficult or easy. The joint interpretation is difficult if we want to rigorously quantify the two-factor or higher-order interaction effects in the functional ANOVA sense since the interaction effect varies locally from one region to another. On the other hand, the joint interpretation is easy if we investigate the feature effects of each LLM separately, which actually involves only main effects within the activation region. In this paper, we take the easy approach while leaving the difficult one for our future work.

We suggest using the parallel coordinate plot for simultaneous visualization of LLM coefficients from different regions. As a common way of visualizing high-dimensional data, the parallel coordinate plot treats each feature as a vertical axis and align all the features in natural ordering (by default) in parallel to each other. Here, each LLM is treated as a subject with its feature coefficients plotted as a connected line from left to right. Repeat it for multiple LLMs on the same chart, and we obtain the graphical output; see Figures 8–9 for the BostonHouse and FicoHeloc examples, where each colored line represents a different LLM with its feature coefficients<sup>4</sup>. Note that these plots are drawn after excluding the LLMs supported by single-instance or single-class regions, which can be identified if the region contains only the instances of the same response value. Later in Section 4 when we discuss the network diagnostics, we will come back to such single-instance or single-class problems, which is a sign of overfitting with instability.

By the parallel coordinate plot, we may check the distribution of main effects from feature to feature and assess their effect consistency among multiple LLMs. For a particular feature, we may find that

- When most of its coefficients are large and have the same signs, it is implied that this feature is an important effect for the final prediction. Furthermore, this feature has a monotonic increasing effect if all coefficients are positive, and a monotonic decreasing effect if all coefficients are negative;
- When most of its coefficients are large and have both positive and negative signs, it is implied that this feature has inconsistent slopes among local regions, which might be caused by either its own nonlinear effect or the interaction effects with other features;
- When most of its coefficients are small and close to zero, it is implied that this feature is not important and can therefore be excluded from the model.

---

<sup>4</sup>As an option, we may also include the intercept of each LLM as the starting vertical axis of the parallel coordinate plot.

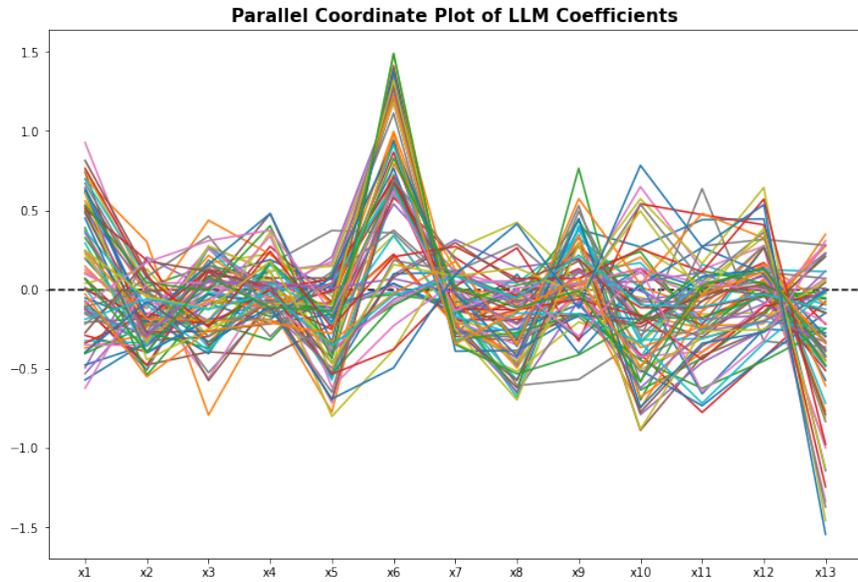


Figure 8: Parallel coordinate plot of unwrapped LLMs excluding dummy cases (Data: BostonHouse; ReLU Net:  $[20]^*2$ )

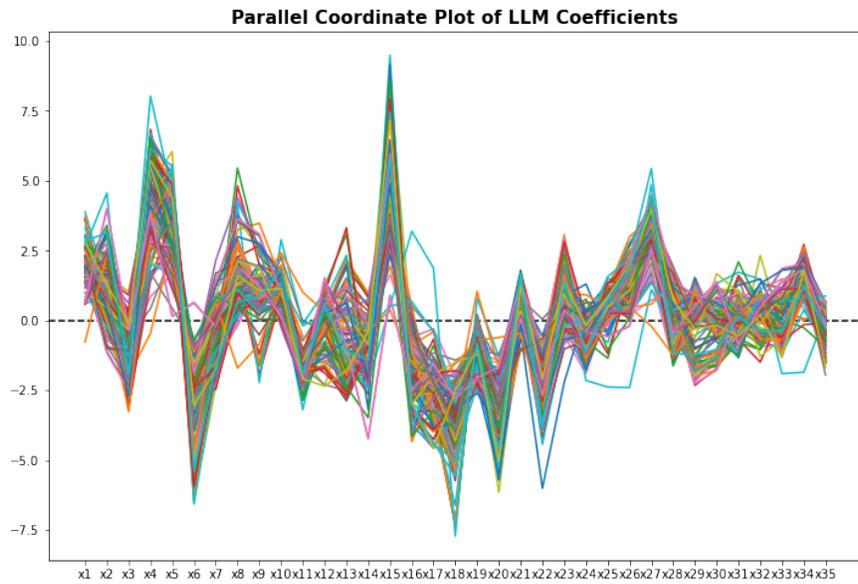


Figure 9: Parallel coordinate plot of unwrapped LLMs excluding dummy cases (Data: Fico-Heloc; ReLU Net:  $[40]^*4$ )

The joint interpretation by parallel coordinate plot is closely related to the feature importance and profile plot in the previous subsection. In an unrealistic special case when all the regions have the same number of training instances, the measure of joint importance (18) coincides with the sum of squared coefficients for each vertical axis in the parallel coordinate plot. Nevertheless, we can check the plots in Figures 8–9 and find that a)  $x_6, x_{13}$  are important features for the BostonHouse example, and b)  $x_{15}, x_4, x_6, x_{18}, x_{27}$  are important features for the FicoHeloc example. This is kind of consistent with the feature importance plot in Figures 6–7 (left). Moreover, if the drill-down analysis for certain features is needed, we can use the local linear profile plot in Figures 6–7 (right).

## 4 LLM-based Network Diagnostics

The unwrapper developed in Section 2 also enables us to perform diagnostics on a pre-trained ReLU DNN based on the set of LLMs. Here we try to address the following three questions:

- Q1. Is it a large or small LLM (in terms of the support size)?
- Q2. Is it a strong or weak LLM (in terms of prediction performance)?
- Q3. Is it a unique or duplicate LLM (in terms of model similarity)?

To answer these questions, for each input of pre-trained ReLU net, the unwrapper in **Aletheia**<sup>©</sup> is designed to generate numerical output that carries the following region-wise summary statistics,

- **Count**: number of training instances within the activation region;
- **Response Mean**: average of response values in the region;
- **Response Std**: standard deviation of response values in the region;
- **Local MSE (regression) or Local AUC (classification)**: performance on the training instances within the region;
- **Global MSE (regression) or Global AUC (classification)**: performance on the entire training instances.

The regions are sorted in the decreasing order of **Count**. See the unwrapper outputs for the CoCircles (classification) and BostonHouse (regression) examples in Figures 10 and 11. Our discussions about network diagnostics are based on such unwrapper outputs.

	<b>Count</b>	<b>Response Mean</b>	<b>Response Std</b>	<b>Local AUC</b>	<b>Global AUC</b>
<b>0</b>	62.0	0.532258	0.498958	0.736677	0.497650
<b>1</b>	60.0	0.600000	0.489898	0.915509	0.496326
<b>2</b>	52.0	0.480769	0.499630	0.828148	0.502093
<b>3</b>	31.0	0.516129	0.499740	0.800000	0.496405
<b>4</b>	29.0	0.655172	0.475312	0.810526	0.504409
...	...	...	...	...	...
<b>525</b>	1.0	1.000000	0.000000	NaN	0.503377
<b>526</b>	1.0	1.000000	0.000000	NaN	0.500440
<b>527</b>	1.0	1.000000	0.000000	NaN	0.499471
<b>528</b>	1.0	1.000000	0.000000	NaN	0.497962
<b>529</b>	1.0	1.000000	0.000000	NaN	0.497462

530 rows × 5 columns

Figure 10: Python code output from **Aletheia**<sup>©</sup> : UnwrapperClassifier (Data: CoCircles)

	<b>Count</b>	<b>Response Mean</b>	<b>Response Std</b>	<b>Local MSE</b>	<b>Global MSE</b>
<b>0</b>	20.0	0.216667	0.056409	0.002481	0.036408
<b>1</b>	17.0	0.465359	0.102063	0.003447	0.098258
<b>2</b>	14.0	0.137778	0.061061	0.002398	0.111933
<b>3</b>	12.0	0.331111	0.052266	0.000367	0.026821
<b>4</b>	12.0	0.167593	0.056426	0.001435	0.031279
...	...	...	...	...	...
<b>150</b>	1.0	0.131111	0.000000	0.000039	0.075982
<b>151</b>	1.0	0.393333	0.000000	0.000116	0.047147
<b>152</b>	1.0	0.428889	0.000000	0.002297	0.044442
<b>153</b>	1.0	0.424444	0.000000	0.000005	0.127390
<b>154</b>	1.0	0.364444	0.000000	0.000654	0.067349

155 rows × 5 columns

Figure 11: Python code output from **Aletheia**<sup>©</sup> : UnwrapperRegressor (Data: Boston-House)

### 4.1 Polar Coordinate Plot

To answer Q1 and Q3, we propose a polar coordinate plot based on the region-wise count statistic and the corresponding LLM coefficients. For two-dimensional data, the bivariate coefficients upon unit norm scaling would determine the direction, while **Count** determines the radius. For multi-dimensional LLM coefficients, we use the principal component analysis (PCA) for dimension reduction to PC1 and PC2, then draw the polar coordinate plot. Besides, the blue color is used to mark the regions with **Response Std** = 0, which corresponds to single-instance regions (regression) or single-class regions (classification). See the bottom part of Figures 10 and 11 about these small regions. In the polar coordinate plot, these small blue-colored regions usually concentrate around the center of plot.

See Figure 12 for the polar coordinate plots for the CoCircles, BostonHouse, and FicoHeloc examples. It is found that the blue-colored regions occupy 84.9%, 52.9%, and 98.2%, respectively. The proportions are very high, especially the two classification examples that involve a great amount of single-class regions. As a matter of fact, the local interpretability discussed in Section 3 is only applicable to large or medium regions exceeding the minimum sample size requirement, but not small regions with few instances.

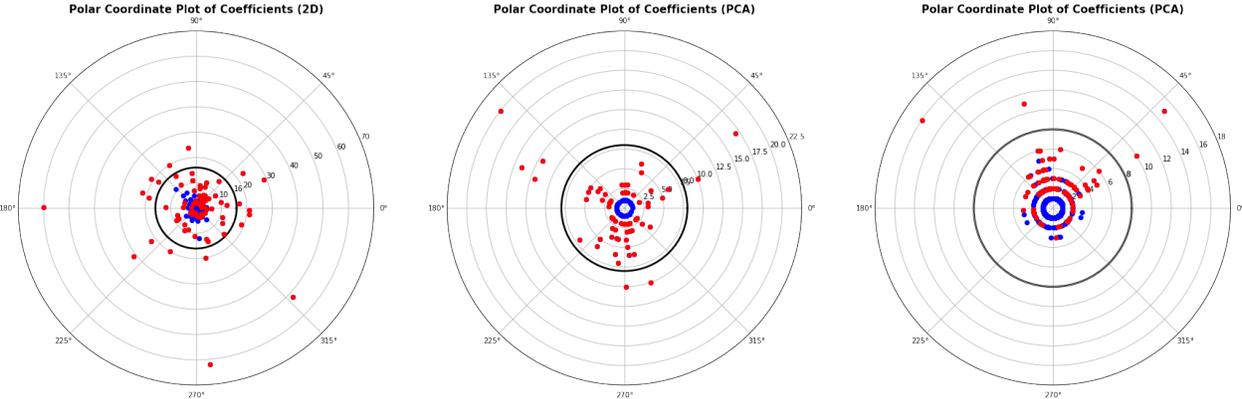


Figure 12: Polar coordinate plot with direction determined by LLM coefficients, while the radius determined by the count of instances. The single-instance or single-class regions are marked in blue. Left: CoCircles; Center: BostonHouse; Right: FicoHeloc.

The polar coordinate plot is also helpful to check the similarity among large/medium LLMs. When multiple LLMs lie in approximately the same direction, it is implied that they share similar model coefficients up to a scaling factor. These similar models could be clustered for network simplification, as we will discuss in Section 5. For the BostonHouse example, along around 30° (also 150°, 270°) direction there are multiple LLMs that are similar to each other. However, for the CoCircles example, the LLMs are seen to scatter around 360 degrees; see also Figure 5 (left).

Those single-instance or single-class regions are resulted from the overparameterized networks, and they are usually redundant. The reliability of the estimated LLMs for these small regions are questionable. Such small regions/LLMs can be either skipped or merged into their neighboring medium/large regions. In Section 5, we discuss a novel merging strategy for dealing with the small or similar LLMs, which would not only simplify the complex network but also lead to enhanced model interpretability and sometimes improved prediction performance.

## 4.2 Local vs. Global Performance

To answer Q2 whether a target LLM is strong or weak, we assess both the local performance within its own region, but also the global performance across other regions. This latter can be referred to *extrapolation strength diagnostics*, which can be done by comparing the local vs. global performances in terms of MSE (regression) or AUC (classification); refer to the last two columns in Figures 10 and 11. Alternatively, we can use the side-by-side bar charts for comparing local and global performances of the top-ranked LLMs; see Figure 13 for the CoCircles, BostonHouse, and FicoHeloc examples.

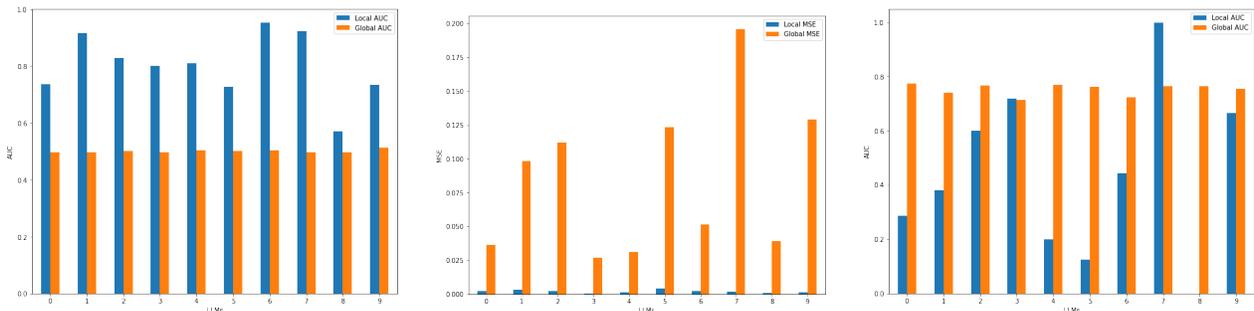


Figure 13: Diagnostics of extrapolation strengths of top-10 LLMs by comparing their local and global performances. Left: CoCircles (AUC); Center: BostonHouse (MSE); Right: FicoHeloc (AUC)

By comparing local vs. global performance of a target LLM, we have the following rules of extrapolation strength diagnostics:

- If performance is good locally but poor globally, it is said to have poor extrapolation strength; e.g. almost all the top-ranked LLMs in the CoCircles and BostonHouse examples in Figure 13.
- If performance is good both locally and globally, it is said to have good extrapolation strength; e.g. the 3, 2, 9-indexed LLMs in the FicoHeloc example in Figure 13.

- If performance is poor locally but good globally, it is said to have extraordinary extrapolation strength; e.g. the 0, 1, 4, 5, 6, 8-indexed LLMs in the FicoHeloc example in Figure 13.

Here for the FicoHeloc example, there exist many LLMs with extraordinary extrapolation strength. It can be found these LLMs share similar model coefficients. Furthermore, as we will discuss in the next section, the input-output relationship in FicoHeloc can be approximated well by a regularized logistic model. In this case, although the ReLU net is highly overparameterized, it can successfully identify the underlying LLMs with a sacrifice on local performance but extraordinary global performance.

## 5 LLM-based Network Simplification

In this section, two strategies are introduced to simplify the extracted LLMs, including merging and flattening. The former reduces the number of LLMs by merging the locally homogeneous LLMs while the latter further simplifies merged LLMs into a single hidden layer ReLU network. For the merged LLMs, we also discuss how to perform local inference with statistical rigor.

### 5.1 Merging

Two LLMs can be merged as their corresponding regions are nearby with similar local linear coefficients and intercepts. To satisfy this principle, agglomerative clustering with connectivity constraint is employed to merge LLMs. In specific, each LLM is initialized as a single cluster and different clusters with similar coefficients and intercepts will be merged hierarchically. The similarity between two LLMs is defined as the Euclidean distance between their local linear coefficients and intercepts, subject to the connectivity constraint. The connectivity matrix is calculated by the  $T$  nearest neighbors for each LLM, according to their region centers

$$\mu_{\mathbf{P}} = \frac{1}{|\mathcal{R}_{\text{train}}^{\mathbf{P}}|} \sum_{\mathbf{x}_i \in \mathcal{R}_{\text{train}}^{\mathbf{P}}} \mathbf{x}_i, \quad \text{for } \mathbf{P} \in \mathcal{P}_{\text{train}}. \quad (19)$$

By incorporating the connectivity constraint in agglomerative clustering, the merged LLMs would satisfy both locality and homogeneity.

After merging, it is possible that some clusters only contain a small number of samples. These small clusters are trivial for interpretation. In practice, we use a post-hoc processing step to merge these small clusters to their nearby large clusters. The small and large clusters can be distinguished by a sample size threshold, e.g.,  $\tau = 30$ . We use `AgglomerativeClustering` in `sklearn.cluster` for agglomerative clustering and the connectivity matrix is obtained by

`kneighbors_graph` in `sklearn.neighbors`. The number of clusters  $K$  can be tuned using grid search, e.g., from  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20\}$ . The number of neighbors  $T$  is initialized to be a small percentage of the total number of LLMs extracted from the training data, e.g. 1% of  $|\mathcal{P}_{\text{train}}|$ . If this value is too small to be satisfied, then we increase it until all LLMs can be eventually connected into one cluster.

Note some new data may not belong to any existing LLMs discovered in the training set, such that the region centers of new LLMs are not available. For simplicity, we use these new samples themselves to represent their corresponding region centers, and the  $T$  nearest LLMs can be accordingly detected.

Finally, we refit a local linear model for each cluster. In specific, two kinds of linear models are considered, i.e., simple GLM for low-dimensional problems or regularized GLM for high-dimensional tasks. The overall merging procedures are summarized in Algorithm 1.

- **Refit model: GLM.** For low-dimensional problems, a simple GLM will be sufficient to capture the local patterns and the resulting model can be largely simplified and is easy to interpret. For the two synthetic datasets, Figures 14 and 15 present the new partitioning regions and local linear models after merging. It can be observed that a lot of the trivial local regions are merged, and the resulting large regions are close to their corresponding ground truth settings.
- **Refit model: Regularized GLM.** High-dimensional tasks often suffer from multicollinearity problems and many coefficients can also be hard to interpret. Therefore, we can use the  $\ell_1$ - or  $\ell_2$ -regularized GLM to make the local linear models more stable and interpretable. In Figure 16, the parallel coordinate plot of merged LLMs is provided for the BostonHouse dataset. According to this figure, 6 clusters are automatically formed. The corresponding plot for the FicoHeloc dataset is provided in Figure 17, in which it is found that only one cluster is left. This means that a deep neural network is over complicated for this dataset, and a logistic regression model is sufficient to capture all the patterns.

## 5.2 Local Inference

The merged LLMs unwrapped from deep ReLU networks can be used for local inference with  $p$ -values and confidence intervals, which leads to an exact and consistent type of local interpretability. In contrast, the state-of-the-art post-hoc interpretation methods including LIME and SHAP [38, 32] are non-exact and varying with input perturbations.

For any input instance  $\mathbf{x} \in \mathcal{X}_{\text{train}}$ , it is straightforward to determine its activation pattern in a pre-trained  $\mathcal{N}_{\text{train}}$ ; so we can locate the local activation region  $\mathbf{x}$  belongs to. The

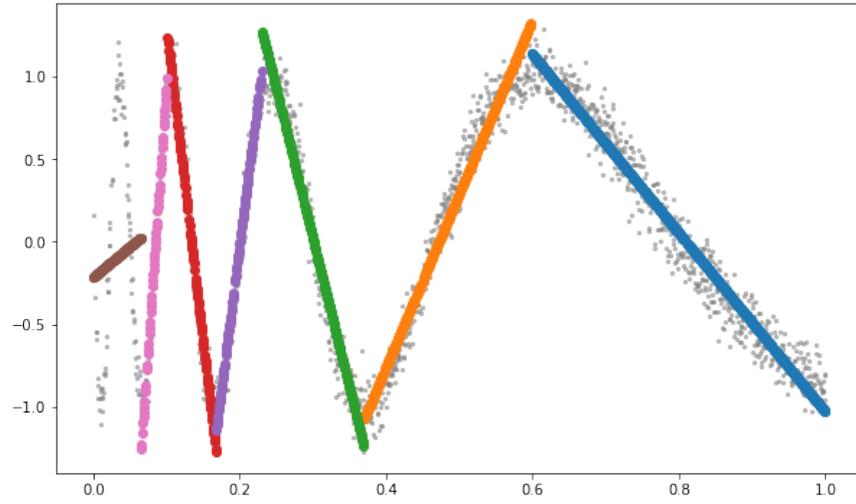
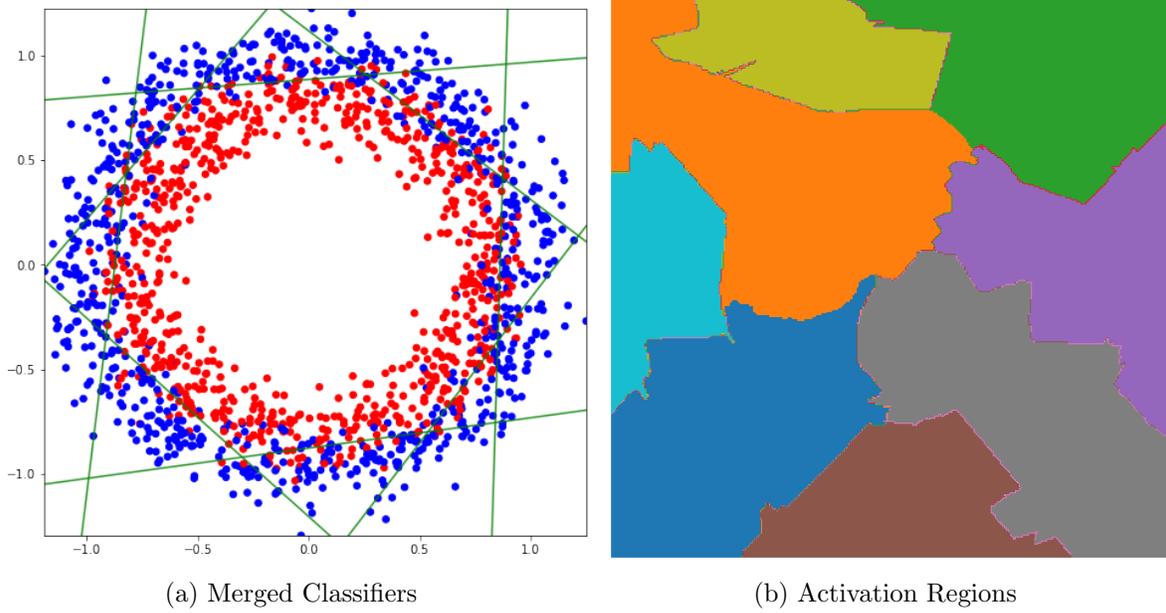


Figure 14: Illustration of ChirpWave after merging.



(a) Merged Classifiers

(b) Activation Regions

Figure 15: Illustration of CoCircles after merging.

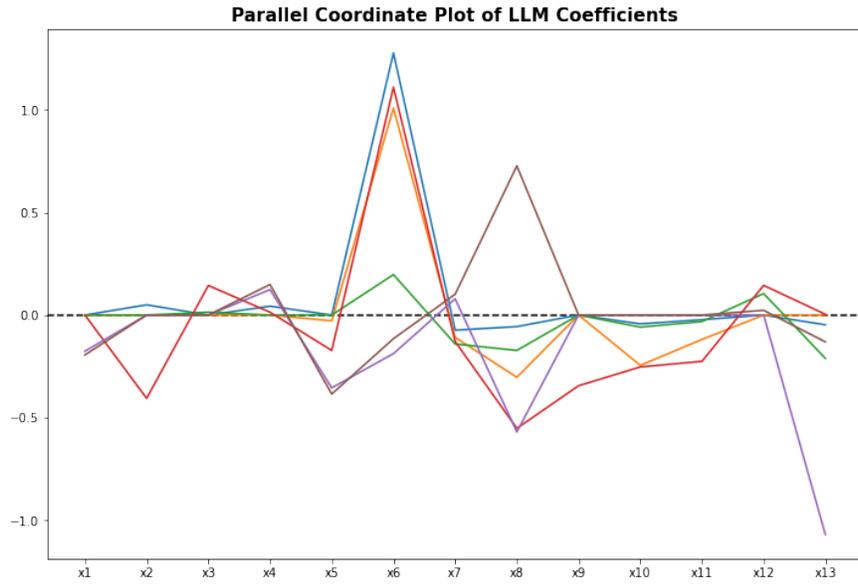


Figure 16: Parallel coordinate plot of the merged network (Data: BostonHouse)

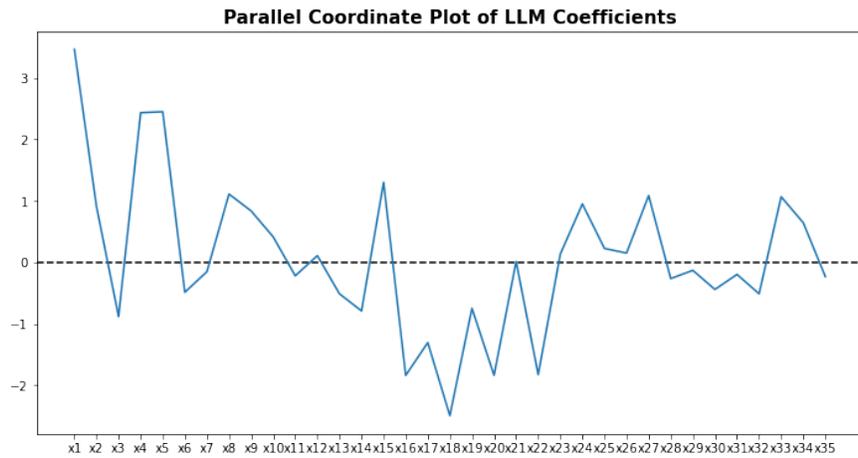


Figure 17: Parallel coordinate plot of the merged network (Data: FicoHeloc)

---

**Algorithm 1** Merging Algorithm

---

**Require:**  $\{\mathbf{x}_i, y_i\}_{i \in [n]}$ , Unwrapped LLMs,  $K$  (Number of Clusters),  $T$  (Number of Neighbors),  $\tau$  (Threshold for Small Clusters).

- 1: Calculate region centers  $\mu_{\mathcal{P}}$  of each LLM.
  - 2: Obtain the  $T$  nearest neighbors for each LLM and construct the connectivity matrix.
  - 3: Cluster  $\{(\tilde{\mathbf{w}}^{\mathcal{P}}, \tilde{b}^{\mathcal{P}})\}_{\mathcal{P} \in \mathcal{P}_{\text{train}}}$  hierarchically into  $K$  groups s.t. connectivity constraint.
  - 4: Find large clusters (sample size larger than  $\tau$ ) and small clusters (the remaining).
  - 5: Merge small LLM clusters to their nearest large clusters.
  - 6: Refit a (regularized) GLM for each new merged LLM.
- 

local interpretability of  $\mathbf{x}$  is translated to its local region with similar instances. When the sample size within the local region is sufficient, statistical inference can be conducted for the corresponding merged LLMs with the local region support.

For a given merged local region, denote the estimated local linear coefficients by  $\hat{\boldsymbol{\beta}} = [\tilde{b}^{\mathcal{P}}, \tilde{\mathbf{w}}^{\mathcal{P}}]$  and the design matrix of  $n_1$  local training instances by  $\mathbf{X} \in \mathbb{R}^{n_1 \times d}$ . For regression problem, the covariance matrix of  $\hat{\boldsymbol{\beta}}$  can be computed by  $\boldsymbol{\Sigma} = \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}$ , in which the variance estimate is given by  $\hat{\sigma}^2 = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 / (n_1 - d - 1)$ . Then, the Wald test of each individual coefficient  $\hat{\beta}_j$  (including the intercept) uses the  $t$ -statistic  $t_j = \hat{\beta}_j / \sqrt{\boldsymbol{\Sigma}_{jj}}$ , which gives  $p$ -value and 95% confidence interval based on the null Student  $t$ -distribution with degrees of freedom  $n_1 - d - 1$ . For binary classification, the local inference can be also conducted following the theory of logistic regression. In this case, the covariance matrix of  $\hat{\boldsymbol{\beta}}$  can be computed by

$$\boldsymbol{\Sigma} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}, \quad \text{with } \mathbf{W} = \text{diag}\left(\{\hat{p}_i(1 - \hat{p}_i)\}_{i=1}^{n_1}\right),$$

where the predicted probability for each instance is given by  $(1 + \exp\{-\hat{\boldsymbol{\beta}}^T \mathbf{x}_i\})^{-1}$ . Then, the Wald test uses the  $z$ -statistic  $z_j = \hat{\beta}_j / \sqrt{\boldsymbol{\Sigma}_{jj}}$  with null distribution  $N(0, 1)$ .

For a quick demonstration of local inference, consider the top two merged LLMs for the CoCircles example, as visualized in Figure 18. For each region, we conducted the local inference procedure with testing results shown in Table 2. The intercept and coefficients are tested all significant (with level 0.05) for both regions. These testing results are reasonable upon a visual check as shown in Figure 18.

Besides, local inference can be performed for  $\ell_1$ -regularized local linear models. For demonstration, the left sub-figure in Figure 19 shows the coefficients and intercepts for the BostonHouse dataset based on 100 bootstrap replicates; while the right sub-figure presents the probability that the coefficient or intercept is reduced to zero. For instance,  $\mathbf{x}_6$  is tested to be a significantly positive feature of the response.

Table 2: Local inference results for top two merged regions (Data: CoCircles).

<b>Region 1</b>	coef	std_err	z	p-value	[0.025	0.975]
intercept	5.3201	0.9076	5.8616	0.0	3.5412	7.0990
x1	4.3638	0.7963	5.4802	0.0	2.8031	5.9245
x2	4.1135	0.8906	4.6188	0.0	2.3679	5.8590
<b>Region 2</b>	coef	std_err	z	p-value	[0.025	0.975]
intercept	5.5461	0.8524	6.5064	0.0000	3.8754	7.2168
x1	4.4345	0.9225	4.8069	0.0000	2.6264	6.2426
x2	-3.4950	0.9619	-3.6335	0.0003	-5.3802	-1.6098

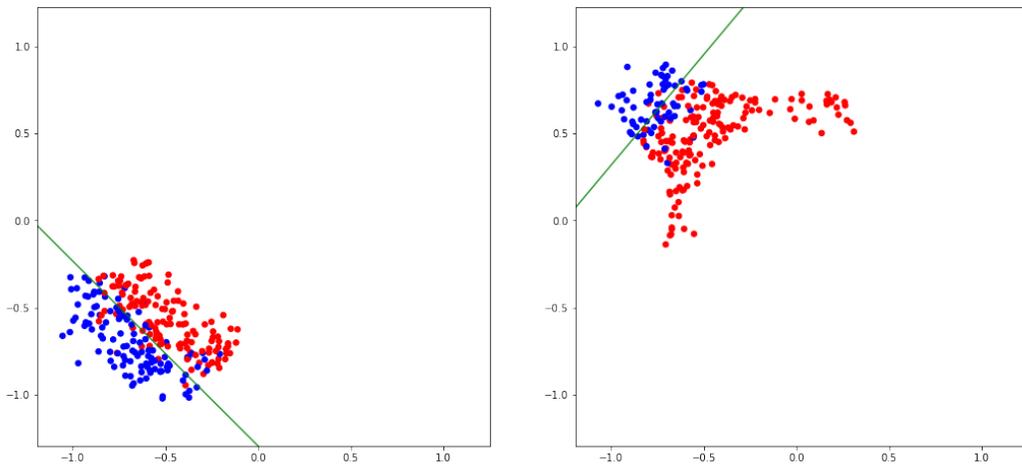


Figure 18: Top two regions of merged LLMs with training instances (Data: CoCircles).

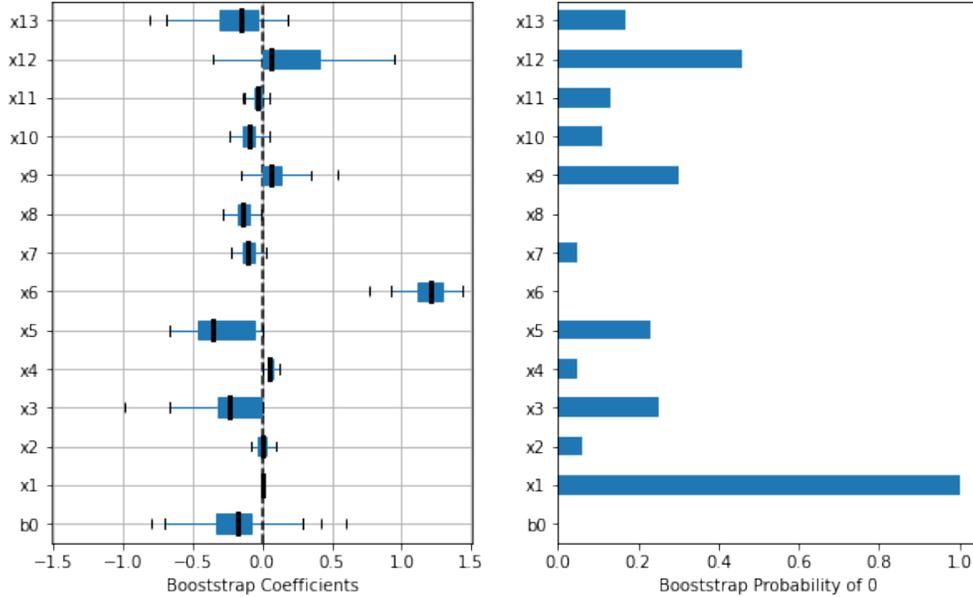


Figure 19: Bootstrap inference for BostonHouse dataset.

### 5.3 Flattening

As LLMs are merged into a small number of clusters, we may do further simplification by flattening them into a single hidden layer ReLU network. That is, a new single hidden layer ReLU network is created where the number of hidden neurons is set to the number of merged LLM clusters. The hidden layer weights and biases are configured by the local linear coefficients and intercepts of the merged LLM clusters. The output layer weights and bias can be estimated via GLM. After that, we fine-tune the flattened network for some epochs, to achieve better predictive performance. In the original deep ReLU network, each LLM corresponds to a linear model that is only active locally. But as the network is flattened, each of the original local regions is now affected by an ensemble of multiple linear models. Below, we present Algorithm 2 for the flattening strategy.

---

**Algorithm 2** Flattening Algorithm

---

**Require:**  $\{\mathbf{x}_i, y_i\}_{i \in [n]}$ , Merged LLM Clusters.

- 1: Collect local linear coefficients and intercepts from each merged LLM cluster.
  - 2: Configure the first hidden layer weights and biases by the extracted local linear coefficients and intercepts, respectively.
  - 3: Forward propagate the data and configure the output layer by GLM.
  - 4: Fine-tune the flattened network using SGD-based optimizers.
- 

In Table 3, we compare the test set predictive performance of the original deep ReLU

Table 3: Average results of ReLU Net, Merged Net, Flattened Net and SLFN, based on 100 Monte Carlo replicates (with standard deviations in parentheses). The last column indicates the average number of clusters used in Merge-Net and FL-Net, as well as the number of hidden nodes in SLFN.

Data	ReLU-Net	Merge-Net	FL-Net	SLFN	n_cluster
ChirpWave (MSE)	0.0486 (0.0356)	0.0721 (0.0540)	0.0839 (0.0323)	0.3351 (0.1021)	10.2800 (2.4498)
CoCircles (AUC)	0.9174 (0.0157)	0.8720 (0.0288)	0.8845 (0.0296)	0.8280 (0.1387)	9.3800 (2.9691)
BostonHouse (MSE)	0.0074 (0.0027)	0.0079 (0.0021)	0.0087 (0.0031)	0.0132 (0.0087)	3.9500 (1.2913)
FicoHeloc (AUC)	0.8002 (0.0091)	0.8050 (0.0088)	0.7962 (0.0107)	0.7545 (0.1102)	1.5500 (1.0897)

network (ReLU-Net), merged LLMs (Merge-Net), flattened network (FL-Net), and single hidden layer feed-forward network (SLFN). Note that SLFN is configured to have the same number of hidden nodes as is used in FL-Net. The only difference between FL-Net and SLFN lies in that they have different initialization methods. That is, the hidden layer of FL-Net is initialized by the refitted local linear models of Merge-Net, while SLFN’s hidden layer is randomly initialized. Compared with the original ReLU-Net, the proposed Merging strategy not only reduces the model complexity but also may improve the generalization performance, e.g., on the FICO dataset. For the other datasets, its test performance is still close to that of ReLU-Net. The flattening strategy can further reduce model complexity. Given the same model complexity, FL-Net has significantly superior predictive performance as compared to SLFN. This indicates the local linear coefficients of Merge-Net make a good initialization for FL-Net.

## 6 Case Study on Home Lending Dataset

In this section, we present a real-world case study of predicting credit default for home lending. The purpose of the case study is to illustrate how the methodology described in this paper is applied in real-life environments. In a regulated financial institution, model interpretability is a requirement for the purpose of evaluating the conceptual soundness of models as well as to explain to model users or customers the decision made by models. Part

of the conceptual soundness requirement, the effects of variables must be consistent with the expected business knowledge. For example, lower credit quality (e.g. FICO score) drives a higher probability of default. The lack of ability to evaluate the conceptual soundness of DNNs has been a major obstacle for credit risk applications. In this case study, we demonstrate that by extracting the equivalent local linear models for ReLU-Net we can provide a strong mechanism to evaluate model conceptual soundness equivalent to the traditional statistical linear regression models. The model diagnostics are applied to identify local linear models that can be improved. Further, model simplification is applied to enhance both model interpretability and model performance. It will be shown below how the large number of local linear models from ReLU-Net are reduced into a smaller number of local linear models and at the same time improve the performance from the original networks.

We modified the dataset in this case study for confidentiality purposes without affecting the ability to demonstrate the methodology. The response variable of interest is binary with the label 0 as “customer does not default” and 1 as “customer defaults” based on various loan characteristics as well as other auxiliary variables totaling 55 predictors where the top 19 are summarized in Table 4. The dataset has 1 million instances, of which about 1% are defaulting cases. We sampled and split the data into training and testing sets, each with approximately  $\frac{1}{3}$  of the total instances, using the same split as explained elsewhere. Furthermore, we make the dataset balanced, by randomly selecting the data from label 0 class to make sure each class has equal distribution in the final dataset.

## 6.1 Model Details

For illustration purposes, we use a feed-forward 2 layer neural with ReLU as the nonlinear activation with 5 neurons in each layer. The model has a performance on the test set of 75.77% accuracy and ROC AUC of 84.29% that are comparable to the performance of best models reported earlier [12]. Table 5 shows the summary of the largest regions such as the number of samples and the performance of local linear models from ReLU-Net. There are a total of 62 regions of which 14 (23%) are single sample instances or a single class.

Excluding the single instances and class, the coefficients of local linear models are shown as a parallel coordinate plot in Figure 20. The top 10 most important variables and the distribution of their coefficients values are shown in Figure 21.

The profile of the two most important variables, i.e., FICO and Forecasted Loan to Value (LTV) are shown in Figure 22. The higher the FICO the less probability of default while the higher the LTV the higher probability of default, consistent with business intuition.

Notice in Table 5 there are regions where the response means are either closer to 0 (non-default), 1 (default), or 0.5 (mix of default and non-default) with varying AUCs locally. To further enhance the model, we apply region merging particularly to address less reliable

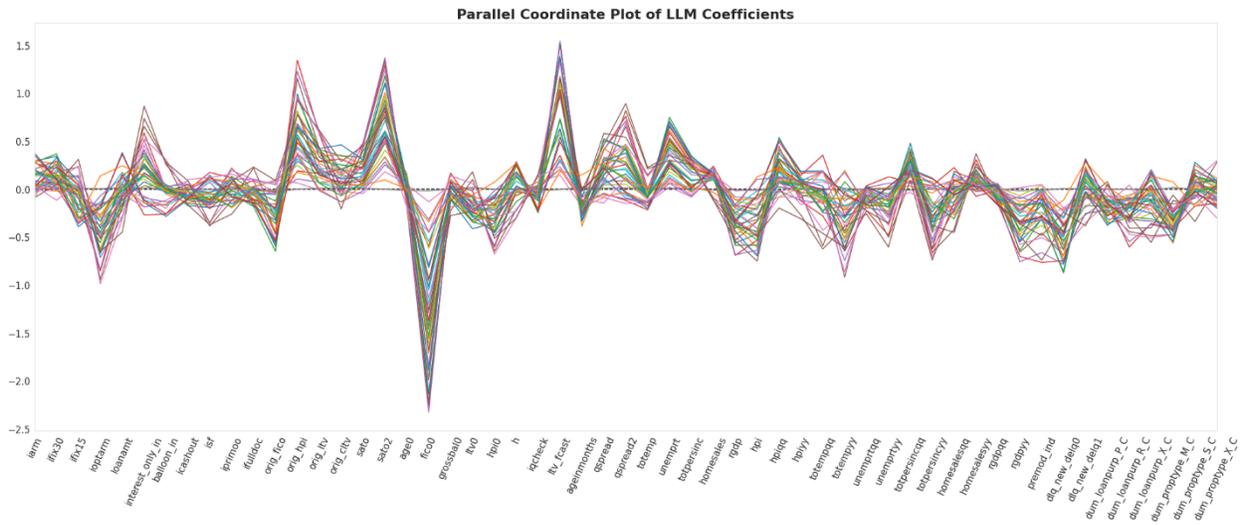


Figure 20: PC plot for the Homelending data.

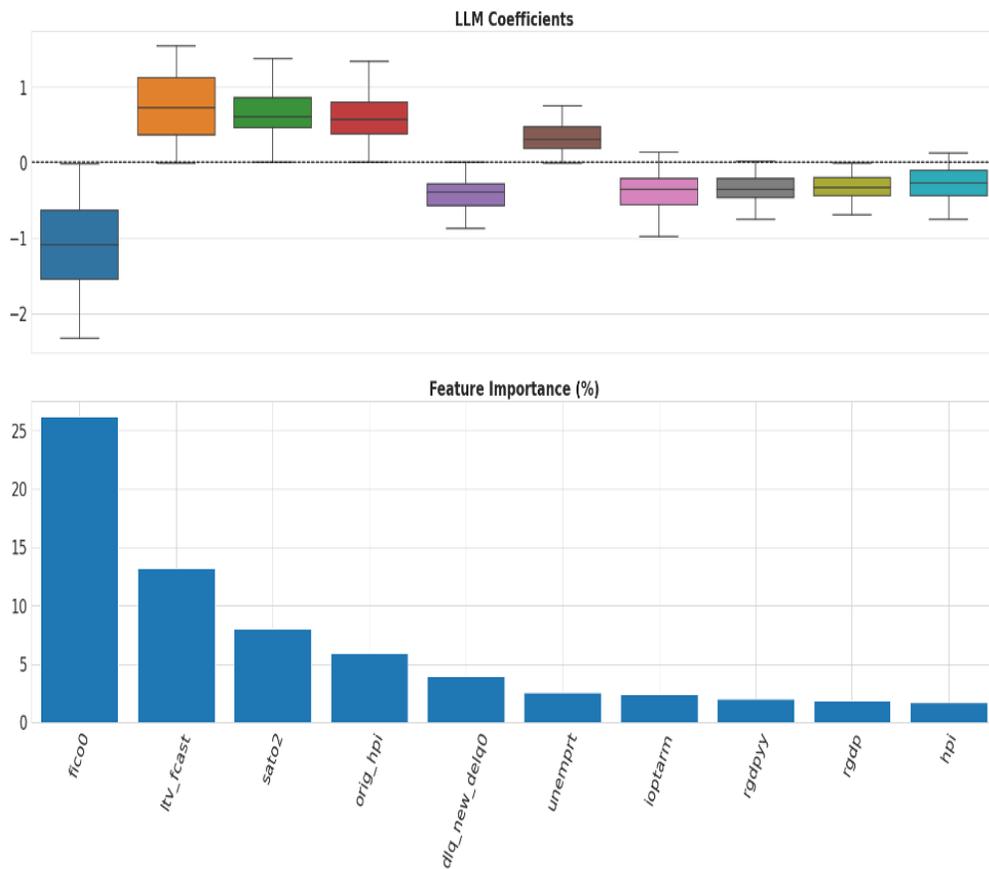


Figure 21: Homelending case study: distribution of LLM coefficients (above); most important features (below).

Table 4: Important features and their definition in the Homelending Dataset.

Variable	Meaning
fico0	FICO at prediction time
ltv_fcast	Loan to value ratio forecasted
sato2	Spread at time of origination
unemprrt	Unemployment rate
dlq_new_delq0	Delinquency status: 1 means current, 0 means delinquent
orig_hpi	Origination house price index
ifix15	Indicator for Fixed 15yr Loan
h	Prediction horizon
premod_ind	Time indicator: before 2007Q2 financial crisis vs. after
qspread	Spread: note rate – mortgage rate
iarm	Indicator for ARM Loan
totpersincyy	Total personal income year over year growth
orig_fico	FICO at origination
hpiyy	HPI YoY
rgdpqq	Real GDP QoQ
orig_cltv	Origination combined ltv
ifulldoc:	Ifulldoc: Indicator for Full Doc
balloon_in	Balloon loan indicator
dum_loanpurp_P	dummy variable of loanpurp=='P'
interest_only_in	Indicator for Interest Only
isf	Single Family

(small) regions where the sample size is small, particularly those regions with single sample size. Here we apply the merging algorithm discussed in 5. In this section, we present the results from the merging algorithm as mentioned in Section 5. The merging step reduced the ReLU-Net regions into three regions as shown in Table 6.

The total AUCs of the merged regions (training AUC= 0.8532, testing AUC = 0.8388), in fact, are better than the original ReLU-Net (training AUC=0.8476, testing AUC=0.8316). This merging step shows that most of the regions created by ReLU-Net are redundant and can be replaced by three LLM that gives better predictive performance. Here, we applied penalized regressions (`sklearn.linear_model.LogisticRegression` with regularization parameter  $C = 0.1$ ) for each merged region.

Notice the results of the merging step created three distinct data segmentation (LLMs)

Table 5: Summary of 15 Largest Regions Unwrapped from ReLU-Net.

	Count	Response Mean	Response Std	Local AUC	Global AUC
0	2303	0.461	0.498	0.791	0.834
1	615	0.668	0.471	0.766	0.824
2	527	0.557	0.497	0.718	0.832
3	478	0.495	0.499	0.699	0.831
4	320	0.043	0.204	0.641	0.829
5	258	0.732	0.442	0.674	0.806
6	218	0.825	0.379	0.628	0.768
7	206	0.038	0.193	0.761	0.784
8	187	0.658	0.474	0.688	0.804
9	155	0.935	0.245	0.579	0.823
10	138	0.217	0.412	0.631	0.820
11	131	0.145	0.352	0.679	0.791
12	121	0.793	0.404	0.775	0.796
13	116	0.129	0.335	0.606	0.815
14	105	0.647	0.477	0.714	0.812
15	103	0.106	0.308	0.590	0.761

with distinct characteristics. Region 0, the largest region with 85.2% of the data contains mixed of default and non-default cases (response mean near 0.5). Region 1 consists of the majority default cases (90%) where Region 2 consists of the majority non-default cases (85.4%). The comparison between local and global AUCs also indicates the distinction among three regions as the local AUCs are better than the global AUCs. It is interesting to observe the parallel coordinate plot of the three LLMs and the corresponding profile of the most important variables shown in Figures 23 and 24, respectively. The profile plots indicate three distinct borrower characteristics: (1) High credit quality (high FICO, low LTV), (2) Middle credit quality (mid FICO and LTV), and (3) Low credit quality (low FICO, high LTV). The parallel coordinate plot provides contrast in terms of significant variables that drive the default event. FICO and LTV are important factors for the majority of the loans (Region 0). Delinquency status is the most important variable for the low credit quality (Region 1) where loans that are not currently delinquent will have a lower probability of default. Premod\_ind which is the indicator of loan origination before or after the financial crisis is important for the higher credit quality (Region 2). Loans originated post-financial crisis are better performing loans as they are subjected to a higher origination standard. It

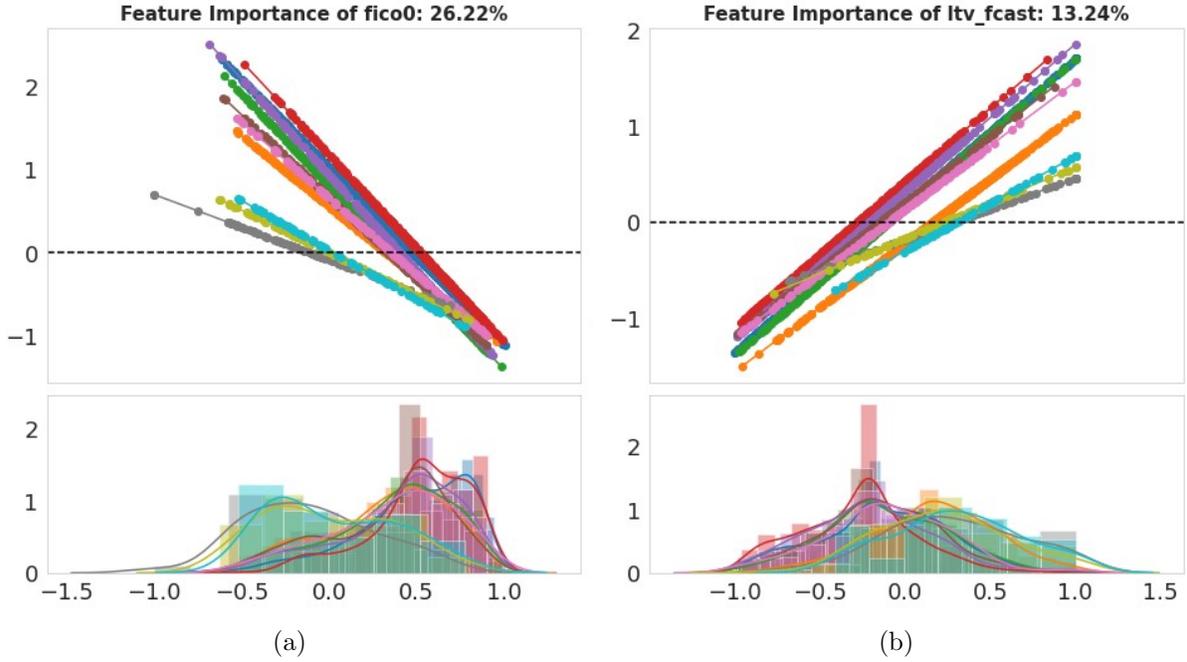


Figure 22: Profile of two most important variables. (a) FICO and (b) Forecasted Loan to Value (LTV).

Table 6: Summary of Three Merged Regions.

Region	Count	Response Mean	Response Std	Local AUC	Global AUC
0	5602	0.460	0.498	0.831	0.854
1	747	0.904	0.293	0.680	0.584
2	219	0.146	0.353	0.987	0.736

is also interesting to point out the effect of the “h” variable which is the time horizon of prediction, an equivalent concept to baseline hazard rate in the survival analysis model. For lower credit quality, the longer time horizon that the loans do not default the less likely the loans will default (i.e., decreasing hazard rate) while the high-quality loan is the opposite which is consistent with business intuition and credit theory from a stochastic process point of view.

In the merging process above we retain the original networks for data partitioning and the final three LLMs are applied through a lookup table to map the regions from the original ReLU-Net to the merged regions. If a smaller network is preferred instead of using the original network, the final LLMs can be used directly as an FL-Net discussed in Section 5.

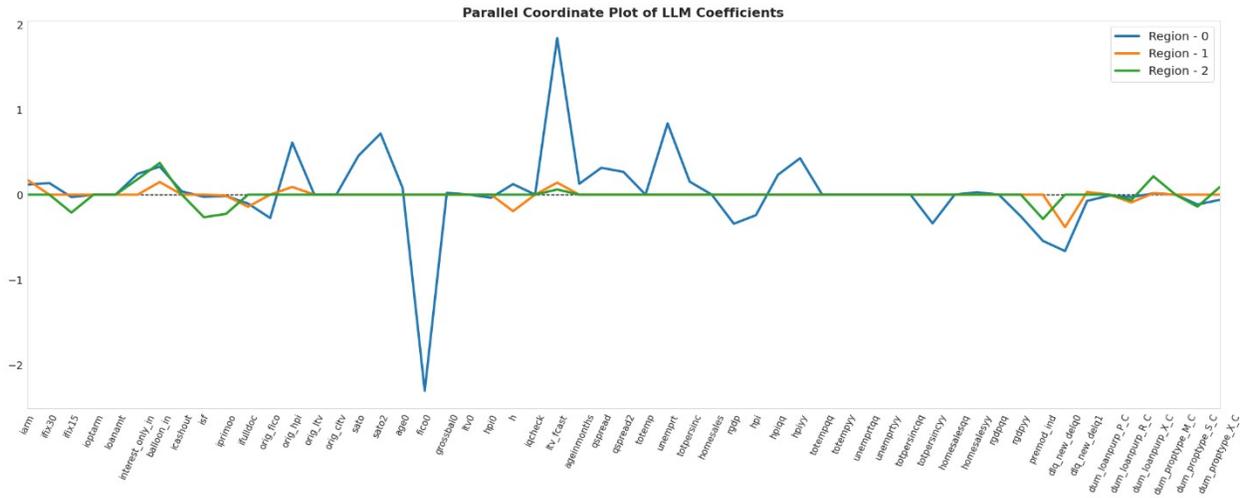


Figure 23: PC plots for the merged regions.

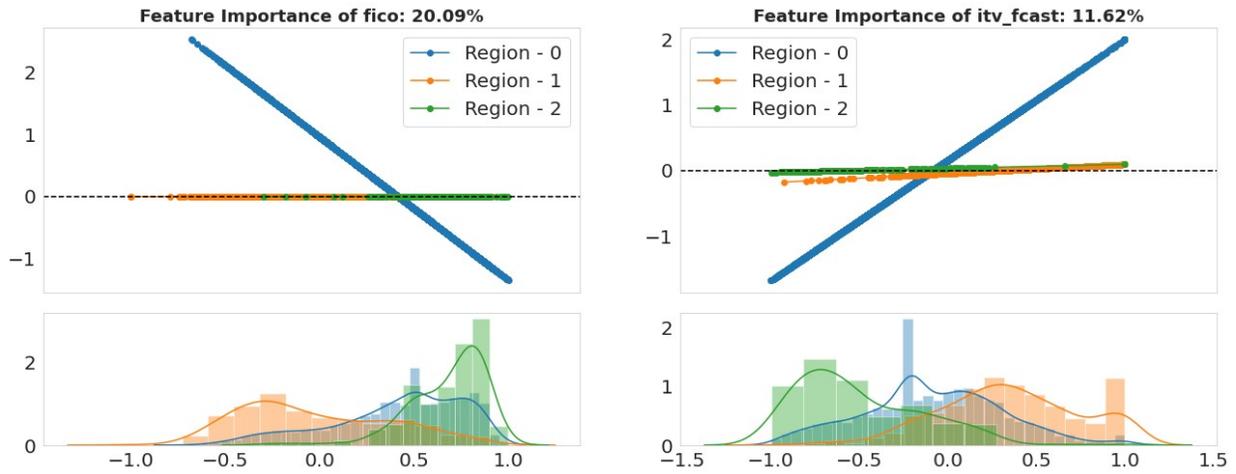


Figure 24: Profile plots of the most important features.

Table 7: AUC Performance Comparison of ReLU-Net, Merge-Net and FL-Net.

	ReLU-Net	Merge-Net	FL-Net
Training	0.8476	0.8532	0.8538
Testing	0.8316	0.8388	0.8368

Here, we constructed a simple single hidden layer network with three nodes to represent the LLMs. The result comparison for the FL-Net is summarized in Table 7. Note that the FL-Net has comparable performance as the Merged-Net and is better than the original ReLU-Net.

## 7 Conclusion

We have introduced a novel **Aletheia**<sup>©</sup> toolkit for unwrapping the black box of deep ReLU networks. Unlike post-hoc explainability methods based on crude assumptions or input permutations, the proposed unwrapper is a rigorous approach to intrinsic interpretability of ReLU DNNs, and it enables us to perform direct interpretation, diagnostics, and simplification based on the unwrapped set of local linear models. Multiple examples used in the paper, including a real case study in home lending credit risk analytics, have demonstrated the novelty and effectiveness of the proposed methods.

## References

- [1] R. Agarwal, N. Frosst, X. Zhang, R. Caruana, and G. E. Hinton. Neural additive models: Interpretable machine learning with neural nets. *arXiv preprint arXiv:2004.13912*, 2020.
- [2] P. Angelov and E. Soares. Towards explainable deep neural networks (xDNN). *Neural Networks*, 130:185–194, 2020.
- [3] D. W. Apley and J. Zhu. Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society Series B*, 82(4):1059–86, 2020.
- [4] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee. Understanding deep neural networks with rectified linear units. In *International Conference on Learning Representations*, 2018.

- [5] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al. Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.
- [6] J. Ba and R. Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pages 2654–2662, 2014.
- [7] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.
- [8] D. Bau, J.-Y. Zhu, H. Strobelt, A. Lapedriza, B. Zhou, and A. Torralba. Understanding the role of individual units in a deep neural network. *Proceedings of the National Academy of Sciences*, 2020.
- [9] M. Belkin, D. Hsu, S. Ma, and S. Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2020.
- [10] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [11] C. Bucilua, R. Caruana, and A. Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541, 2006.
- [12] J. Chen, J. Vaughan, V. Nair, and A. Sudjianto. Adaptive explainable neural networks (AxNNs). *Available at SSRN 3569318*, 2020.
- [13] L. Chu, X. Hu, J. Hu, L. Wang, and J. Pei. Exact and consistent interpretation for piecewise linear neural networks: A closed form solution. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1244–1253, 2018.
- [14] J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018.
- [15] J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [16] J. H. Friedman and W. Stuetzle. Projection pursuit regression. *Journal of the American statistical Association*, 76(376):817–823, 1981.

- [17] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323, 2011.
- [18] A. Goldstein, A. Kapelner, J. Bleich, and E. Pitkin. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics*, 24(1):44–65, 2015.
- [19] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio. *Deep learning*. MIT press Cambridge, 2016.
- [20] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [21] D. Gopinath, H. Converse, C. Pasareanu, and A. Taly. Property inference for deep neural networks. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 797–809. IEEE, 2019.
- [22] D. Gunning and D. W. Aha. DARPA’s explainable artificial intelligence program. *AI Magazine*, 40(2):44–58, 2019.
- [23] B. Hanin and D. Rolnick. Deep ReLU networks have surprisingly few activation patterns. In *Advances in Neural Information Processing Systems*, pages 361–370, 2019.
- [24] T. J. Hastie and R. J. Tibshirani. *Generalized Additive Models*, volume 43. CRC press, 1990.
- [25] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [26] J.-N. Hwang, S.-R. Lay, M. Maechler, R. D. Martin, and J. Schimert. Regression modeling in back-propagation and projection pursuit learning. *IEEE Trans. Neural Networks*, 5(3):342–353, 1994.
- [27] J. D. Janizek, P. Sturmfels, and S.-I. Lee. Explaining explanations: Axiomatic feature interactions for deep networks. *arXiv preprint arXiv:2002.04138*, 2020.
- [28] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks. In *International Conference on Computer Aided Verification*, pages 97–117. Springer, 2017.
- [29] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.

- [30] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view from the width. In *Advances in neural information processing systems*, pages 6231–6239, 2017.
- [31] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee. From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1):56–67, 2020.
- [32] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- [33] P. McCullagh and J. A. Nelder. *Generalized Linear Models*, volume 37. CRC Press, 1989.
- [34] C. Molnar. *Interpretable Machine Learning*. Lulu.com, 2020.
- [35] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio. On the number of linear regions of deep neural networks. In *Advances in neural information processing systems*, pages 2924–2932, 2014.
- [36] W. J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yu. Definitions, methods, and applications in interpretable machine learning. *Proceedings of the National Academy of Sciences*, 116(44):22071–22080, 2019.
- [37] M. Raghu, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pages 2847–2854. PMLR, 2017.
- [38] M. T. Ribeiro, S. Singh, and C. 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. ACM, 2016.
- [39] J. Schmidt-Hieber et al. Nonparametric regression using deep neural networks with ReLU activation function. *Annals of Statistics*, 48(4):1875–1897, 2020.
- [40] T. Serra, A. Kumar, and S. Ramalingam. Lossless compression of deep neural networks. *arXiv preprint arXiv:2001.00218*, 2020.
- [41] T. Serra, C. Tjandraatmadja, and S. Ramalingam. Bounding and counting linear regions of deep neural networks. In *International Conference on Machine Learning*, pages 4558–4566. PMLR, 2018.

- [42] A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating activation differences. In *ICML*, 2017.
- [43] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju. Fooling LIME and SHAP: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186, 2020.
- [44] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *ICML*, 2017.
- [45] M. Tsang, D. Cheng, and Y. Liu. Detecting statistical interactions from neural network weights. *arXiv preprint arXiv:1705.04977*, 2017.
- [46] J. Vaughan, A. Sudjianto, E. Brahim, J. Chen, and V. N. Nair. Explainable neural networks based on additive index models. *The RMA Journal*, pages 40–49, 2018.
- [47] Z. Yang, A. Zhang, and A. Sudjianto. Enhancing explainability of neural networks through architecture constraints. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [48] Z. Yang, A. Zhang, and A. Sudjianto. GAMI-Net: an explainable neural network based on generalized additive models with structured interactions. *arXiv preprint arXiv:2003.07132*, 2020.
- [49] Q. Zhang, Y. N. Wu, and S.-C. Zhu. Interpretable convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8827–8836, 2018.
- [50] D. Zou, Y. Cao, D. Zhou, and Q. Gu. Gradient descent optimizes over-parameterized deep ReLU networks. *Machine Learning*, 109(3):467–492, 2020.