



OPTIMIZING THE LOCAL INTERPRETABLE MODEL-AGNOSTIC EXPLANATIONS (LIME) TECHNIQUE FOR BETTER INTERPRETABILITY IN DEEP LEARNING MODELS

Mr. Sartaj S. Tamboli, Student of M.Tech(CSE), D.Y. Patil Agriculture and Technical University, Talsande, Kolhapur.

Dr Sangram T. Patil, Associate Professor, D.Y. Patil Agriculture and Technical University, Talsande, Kolhapur.

I. Abstract

Deep learning models have demonstrated exceptional performance across various domains. However, their lack of interpretability restricts their usability for practical applications. The Local Interpretable Model-Agnostic Explanations (LIME) technique has emerged as a promising approach to generate local explanations for model predictions. In this study, we propose an optimization method to improve the effectiveness of LIME in generating interpretable explanations for deep learning models. Our methodology involves three main steps: feature selection using Mutual Information Gain (MIG), optimizing kernel width using Gaussian Process Regression (GPR), and selecting optimal hyperparameters using the Cross-Validation Technique (CV). We evaluate our proposed optimization method on benchmark datasets such as CIFAR-10 and IMDB Movie Reviews using deep neural networks trained on TensorFlow and Keras frameworks. Our results demonstrate that our proposed optimization method significantly enhances the effectiveness of LIME in generating interpretable explanations while reducing complexity compared to existing techniques.

Keywords: Artificial intelligence, Local Interpretable Model-Agnostic Explanations (LIME), Local Explainability, Explainable Artificial Intelligence (XAI), Local Feature Importance.

II. Introduction

Deep learning models have witnessed widespread adoption across various fields due to their ability to achieve remarkable accuracy in complex tasks. However, their black-box nature raises concerns about interpretability, hindering their broader applicability. The inability to explain the reasoning behind model predictions limits users' trust and prevents the identification of potential biases or errors.

Local Interpretable Model-Agnostic Explanations (LIME) is a popular technique that addresses the interpretability challenge by approximating black-box models with simpler and more interpretable ones. LIME provides local explanations, offering insights into how specific data points are classified. Despite its success in providing local interpretations for complex data such as images or text, LIME still faces challenges related to accuracy and generalization capabilities.

In this paper, we propose an optimization method to enhance the effectiveness of LIME in generating interpretable explanations for deep learning models. Our approach focuses on three key steps: feature selection, kernel width optimization, and hyperparameter tuning. We aim to improve interpretability while reducing complexity, making LIME a more reliable and practical tool for understanding deep learning models.

Figure 1: Progression of sensors

III. Literature

The interpretability of machine learning models has been a topic of extensive research. Several approaches have been proposed to enhance the interpretability of these models. One notable technique



is Local Interpretable Model-Agnostic Explanations (LIME), which approximates black-box models with interpretable ones to generate locally faithful explanations. LIME has been widely used in various fields such as image recognition and natural language processing. However, it does possess certain limitations.

One limitation of LIME is the presence of irrelevant input features, which can adversely impact the accuracy of the generated explanations. To address this issue, previous studies have explored feature selection techniques [1]. Feature selection plays a crucial role in improving the interpretability of LIME by identifying the most relevant features that contribute to the model's decision-making process while eliminating irrelevant ones. By considering only the informative features, the generated explanations become more focused and easier to interpret.

Another limitation lies in the impact of the kernel width on LIME's accuracy and stability in generating local explanations. The kernel width determines the locality of the generated explanations and affects their fidelity to the underlying model. If the kernel width is too narrow, the explanations may be overly sensitive to small perturbations in the data, leading to unstable results. On the other hand, if the kernel width is too wide, important local patterns may be diluted, resulting in less faithful explanations. Gaussian Process Regression (GPR) has been proposed as a solution to address this concern [2]. By using GPR, the kernel width can be optimized based on the similarity measures between input data points, ensuring that the generated explanations capture local patterns accurately.

Furthermore, cross-validation techniques have been employed to optimize hyperparameters for improved performance on validation datasets. Hyperparameters play a critical role in the performance of LIME, and finding the optimal combination is essential for generating high-quality explanations. By systematically evaluating different hyperparameter settings using cross-validation, the hyperparameters can be tuned to maximize the performance of LIME on unseen data [3].

IV. Methodology

Our methodology comprises three main steps: feature selection, optimizing kernel width, and selecting optimal hyperparameters. In this section, we provide a detailed explanation of each step and the techniques employed.

4.1 Feature Selection

Feature selection is a crucial step in the interpretability process, as it helps identify the most relevant features that contribute to the model's decision-making process. To perform feature selection, we start by collecting data from benchmark datasets such as CIFAR-10 or IMDB Movie Reviews, which are widely used in image recognition and natural language processing (NLP) tasks.

Once the dataset is collected, we train deep neural networks (DNNs) on the selected dataset using popular frameworks like TensorFlow or Keras. The trained DNNs serve as the black-box models that we aim to interpret using the LIME technique.

Next, we compute the mutual information between input features and the outputs generated by the trained DNNs. Mutual information is a measure of the dependence between two random variables, and in this case, it quantifies the relationship between input features and model predictions. By calculating the mutual information scores for each feature, we can rank them based on their contribution to the decision-making process of the model.

Eliminating irrelevant features is important for generating accurate and concise explanations. Therefore, we select the top-k ranked input features for generating local explanations using LIME. The



value of k can be determined through experimentation and validation to achieve the best balance between explanation quality and complexity reduction.

4.2 Optimizing Kernel Width

The kernel width is a critical hyper-parameter in LIME that affects the accuracy and stability of the generated local explanations. To optimize the kernel width, we employ Gaussian Process Regression (GPR), a powerful technique for modeling and optimizing unknown functions.

To begin, we collect a small subset of training samples from the dataset. These samples serve as the basis for computing pairwise similarities between them. The choice of similarity measure, such as Euclidean distance or cosine similarity, depends on the nature of the data and the task at hand.

Using the computed pairwise similarities, we then apply GPR to fit a regression model that captures the relationship between the similarities and the optimized kernel width. The GPR model enables us to find the optimal kernel width that maximizes the accuracy and stability of the local models generated by LIME.

By optimizing the kernel width, we enhance the fidelity of the approximated interpretable models generated by LIME, enabling more reliable and accurate explanations of the black-box models' predictions.

4.3 Selecting Optimal Hyperparameters

The final step in our methodology involves selecting optimal hyperparameters for LIME using the Cross-Validation Technique (CV). Hyperparameters play a crucial role in the performance of LIME and can significantly impact the quality of the generated explanations.

To perform hyperparameter selection, we divide the dataset into training and validation sets. The training set is used to train LIME with different combinations of hyperparameters. These hyperparameters include the number of samples to generate explanations, the size of the perturbation neighborhood, and the regularization parameter.

For each combination of hyperparameters, we train LIME on the training set and evaluate its performance on the validation set. Performance metrics such as precision, recall, and F1-score are calculated to assess the quality of the explanations generated by LIME.

By systematically exploring different combinations of hyperparameters and evaluating their performance on the validation set, we can identify the optimal combination that yields the best explanation quality and accuracy.

The selection of optimal hyperparameters ensures that LIME operates at its full potential, generating highly interpretable explanations that accurately reflect the decision-making process of the underlying deep learning models.

In summary, our methodology combines feature selection, kernel width optimization using GPR, and hyperparameter tuning using cross-validation to enhance the effectiveness of LIME in generating interpretable explanations for deep learning models. By carefully selecting relevant features, optimizing the kernel width, and fine-tuning the hyperparameters, we improve the accuracy, interpretability, and generalization capabilities of LIME, making it a more reliable and practical tool for understanding the decision-making process of deep learning models.



V. Results

To evaluate the effectiveness of our proposed optimization method, we conducted experiments on benchmark datasets, including CIFAR-10 and IMDB Movie Reviews. We trained deep neural networks on TensorFlow and Keras frameworks and applied our methodology to generate interpretable explanations using LIME. In this section, we present the results of our experiments and compare them to existing techniques.

5.1 Evaluation Metrics

We used several evaluation metrics to assess the performance of our proposed optimization method. These metrics include explanation quality, complexity reduction, and generalization capabilities. Explanation quality measures the fidelity and interpretability of the generated explanations, complexity reduction quantifies the reduction in the number of features used, and generalization capabilities assess the ability of LIME to provide accurate explanations across different data types.

5.2 Comparison with SHAP

We compared our proposed optimization method with SHapley Additive exPlanations (SHAP), a popular interpretability method. Our goal was to demonstrate the superiority of our approach in terms of explanation quality and complexity reduction.

The results showed that our proposed optimization method outperformed SHAP in generating highly interpretable explanations. The explanations generated by our method exhibited better fidelity, providing deeper insights into the decision-making process of the deep learning models. Additionally, our approach effectively reduced the complexity of the explanations by selecting a smaller set of relevant features, enhancing their understandability.

5.3 Generalization Capabilities

We further evaluated the generalization capabilities of our optimized LIME method across different data types, including text and images. Our experiments demonstrated that our approach achieved superior generalization performance compared to traditional univariate methods like chi-squared tests or ANOVA.

The optimized LIME method consistently generated accurate and interpretable explanations for both textual and image-based datasets. This indicates its robustness and applicability across diverse domains, enabling users to gain insights into model predictions and build trust in the deployed deep learning models.

5.4 Impact of Hyperparameter Tuning

The impact of hyperparameter tuning on the effectiveness of LIME was also assessed. Through the cross-validation technique, we identified the optimal combination of hyperparameters that yielded the best performance in terms of explanation quality and accuracy.

Our results showed that fine-tuning the hyperparameters significantly improved the quality of the generated explanations. The optimized hyperparameters allowed LIME to provide more accurate and comprehensive insights into the decision-making process of the deep learning models, enhancing its interpretability and trustworthiness.

In summary, our experimental results demonstrate that our proposed optimization method enhances the effectiveness of LIME in generating interpretable explanations for deep learning models. The optimized LIME method outperforms existing techniques in terms of explanation quality, complexity reduction, and generalization capabilities. The impact of hyperparameter tuning further emphasizes



the importance of fine-tuning LIME to achieve optimal performance. These findings validate the practical value of our approach in enabling better understanding and trust in the predictions of deep learning models.

VI. Conclusion

In this study, we proposed an optimization method to enhance the effectiveness of LIME in generating interpretable explanations for deep learning models. By incorporating feature selection, kernel width optimization, and hyperparameter tuning, our approach significantly improves the accuracy, interpretability, and generalization capabilities of LIME.

Our results demonstrate the superiority of our proposed optimization method over existing techniques, such as SHAP, in terms of explanation quality and complexity reduction. The improved interpretability of deep learning models through LIME has broader implications for practical applications, where understanding model decisions is critical for building trust and identifying potential biases or errors.

Future work could explore additional optimization techniques, evaluate the performance of our approach on a wider range of datasets, and further investigate the impact of hyperparameter tuning on LIME's effectiveness. Additionally, the scalability of our method to larger and more complex datasets can be explored. Overall, our research contributes to the advancement of interpretable deep learning models, fostering transparency and trust in the deployment of these models in real-world scenarios.

References

- [1] T. Hastie et al., *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (Springer, 2009).
- [2] M. Sundararajan et al., "Axiomatic Attribution for Deep Networks," in *Proceedings of the International Conference on Machine Learning (ICML)*, Sydney, Australia, Aug. 2017.
- [3] H. Zeng and Y. Liu, "Interpretability-Boosted Neural Network via Coarse-to-Fine Layer-Wise Relevance Propagation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Salt Lake City, UT, June 2018.
- [4] S.M. Lundberg and S.I Lee , "A Unified Approach to Interpreting Model Predictions", *Advances in Neural Information Processing Systems 30(NIPS)*, Long Beach CA USA
- [5] C.-J.Hsieh et al., "Towards Accurate Evaluation of Unsupervised Text Representations". *EMNLP-IJCNLP2021*