

Want robust explanations? Get smoother predictions first.

Deddy Jobson

Mercari Inc., Roppongi Hills Mori Tower, 6 Chome-10-1 Roppongi, Minato City, Tokyo 106-6118

Abstract

Model-agnostic machine learning interpretability methods like LIME which explain the predictions of elaborate machine learning models suffer from a lack of robustness in the explanations they provide. Small targeted changes to the input can result in large changes in explanations even when there are no significant changes in the predictions made by the machine learning model. This is a serious problem as it undermines the trust one has in the explanations made. We propose to solve the problem by smoothing the predictions of the machine learning model as a preprocessing step. We smoothen the predictions by taking multiple samples from the neighbourhood of each input data point and averaging the output predictions. Through our preliminary experiments, we show that the explanations are more robust because of smoothening thus making them more reliable.

Keywords

interpretable machine learning, model agnostic, interpretability, LIME, robustness

1. Introduction

The sudden improvement in performance of machine learning through deep learning and tree ensemble methods has led to an explosion in the adoption of machine learning in a wide variety of prediction tasks in multiple domains like image, text, tabular data, etc. While the increased performance has made machine learning models much more useful in practice, it has come at the cost of interpretability; one can no longer trivially explain the decisions made by machine learning models the same way one could for statistical models like linear regression in the past. While we can do without interpretability in cases where the consequences of the downstream decisions are little, like in the case of recommending movies, interpretability becomes important in high-stakes situations like predicting whether or not a person has cancer[1]. In such a case, it is not just important to know what the predictions of the model are, but also how the predictions were made.

A number of model-agnostic interpretability methods exist to help explain the predictions made by machine learning models. Partial Dependence Plots[2] show the marginal effect of a feature on the outcome. Individual Conditional Expectation plots[3] do the same by making separate plots for each individual thus allowing one to see the variance (and not just the mean) of the effect of each feature. The above two have a problem wherein we consider the effect of very unlikely counterfactual scenarios in the case where the features in the dataset are strongly correlated.

Shapley values[4] take a game-theoretic approach and assume different features take part in a collaboration to assign a score for an instance. The shapley value for a feature is the average increment in the score obtained by the inclusion of said feature in the collaboration. While using shapley values has a strong mathematical foundation, it has the downside where the computational cost for calculation is exponential to the number of features. While methods like Tree SHAP[5] exist to more efficiently calculate the values, there are issues with the robustness[6] of shapley values which have not yet been resolved.

Local Interpretable Model-Agnostic Explanations (LIME)[7] is a method that estimates a local surrogate model in the vicinity of each data point and uses the coefficients of the local model to interpret the decisions made by the model. It is related to SHAP through Kernel SHAP[8], a way to get approximate SHAP values. One advantage of LIME over shapley values is that LIME can produce sparse explanations which don't rely on too many features resulting in more human-friendly explanations. However, issues regarding the robustness[6] of the explanations provided by LIME have been raised. Our goal in this paper is to find ways to improve the robustness of the interpretations made by LIME to improve the reliability and therefore trustworthiness of the provided explanations.

2. Problem Setup

The original LIME algorithm works as follows, given a trained model and a target data point:

1. Sample data around the neighbourhood of the data point.
2. Get the predicted values for the sampled data points.

AIMLAI '22: Advances in Interpretable Machine Learning and Artificial Intelligence, October 21, 2022, Atlanta, GA

✉ deddy@mercari.com (D. Jobson)

🆔 0000-0003-1557-8131 (D. Jobson)

© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

3. Fit a surrogate model to the generated data weighted by distance from the target data point.
4. Explain the prediction of the main model with the coefficients of the surrogate model.

The explanations generated by the above algorithm can be unstable for a number of reasons.

One source of instability is the sampling of data points[9] that is done randomly, ignoring any correlation between features. Methods have been developed to estimate the required number of samples to get stable explanations[10] or do away with randomness in the sampling altogether[11].

Another potential cause for instability in explanations, especially pertinent to the case of tabular data, is the discretization of the numerical features. While for the most part this can yield more consistent explanations, target data points near the boundaries can have unstable explanations even when the model predictions (which don't rely on discretization) in the vicinity are relatively stable.

3. Related Work

The measurement of the stability (or lack thereof) of LIME's explanations isn't a new research problem. Alvarez-melis et al.[6] have shown that small perturbations to the input can cause a large change in the output without much of a change in the predictions made by the model. They use the definition of Lipschitz continuity to get the maximum possible difference in explanation within the neighbourhood of the data point to be explained. Their approach is similar to prior work that was done to inspect the lack of robustness of predictions made by neural networks[12].

Visani et al.[13] introduce two novel metrics grounded in statistics to measure the extent to which repeated sampling of the data leads to a variance in the explanations. Their metrics quantify the variance of the selected features and coefficient values, the lower the better.

Much more recently, Garreau et al.[14] performed a very deep analysis into the workings of LIME for tabular data and (among other things) found that when the surrogate model (the one trained for interpretability) uses ordinary least squares, and the number of sampled data points is large, the estimations by LIME are robust to mild perturbations. This suggests that the cause of instability could lie elsewhere.

4. Our Method

For our method, we smoothen the predictions of the model we want to explain with the help of Gaussian noise. We do so because we hypothesize that the lack

Table 1

Preliminary experiments on the Boston dataset (the lower the score the better)

Algorithm	Lipschitz Discontinuity Score
LIME	2.78
LIME smoothed	2.60

of robustness in the explanations caused by LIME is not because of LIME itself but rather the jaggedness of the predictions made by the model.

We smoothen the predictions by averaging the predictions made on random perturbations on the data points. We consider the case where all features of the data point are numeric and continuous in this study. We perturb each feature by adding it with gaussian noise of zero mean. We refer to the standard deviation of the gaussian noise to be the "strength" parameter. This is because the greater the "strength" parameter, the larger the perturbations and the smoother the averaged predictions will be (assuming enough samples) and so the "stronger" the smoothening effect. We choose a strength value of 0.1 for our experiments and take 100 random samples for each data point for the smoothening process.

5. Experiments and Discussion

Our hypothesis is that smoothening the predictions will yield explanations that are more robust. To test this hypothesis, we look at the extent to which the variance of LIME's explanations change before and after smoothening the predicting function. We define a metric called Lipschitz Discontinuity Score (LDS) Score which is derived from the expression used in the definition of Lipschitz Continuity. Our approach is similar to the one used in [6]. LDS is defined as follows:

$$LDS = \frac{1}{N} \sum_{i=1}^N \max_{j \neq i} \frac{\|f(x_i) - f(x_j)\|_2}{\|x_i - x_j\|_2} \quad (1)$$

In the above expression, N is the number of records in the dataset, i and j are indices to denote individual records and take values from 1 to N, and $f(x_i)$ is the vector of coefficients we get from the explanations of the LIME algorithm.

We perform preliminary experiments on the publicly available Boston dataset, a dataset with 12 covariates for a regression problem. We parameterize the LIME algorithm to explain with only 3 features. The base model used is the random forest regressor from scikit-learn. We use the default parameters of the random forest since it suffices for the purposes of this study. We estimate the LDS on the Boston dataset using 10-fold cross validation. In table 1, we compare the LDS of the explanations of LIME

for two cases: with and without smoothening. We find that there is a substantial improvement in the LDS when smoothening the predictions, in line with our hypothesis.

6. Future Work

In this paper, we smoothen the predictions of the machine learning model by sampling neighbouring points randomly multiple times and taking the average of the output. We do this to increase the robustness of the explanations by LIME. We chose white noise since the approach is similar to the original LIME algorithm, but since its introduction, various improved sampling strategies have been proposed that result in more robust explanations[15, 16]. Trying those other sampling methods for the purpose of smoothening the predictions is beyond the scope of this extended abstract and can be considered as one avenue for future research.

While we perform preliminary experiments with tabular data, our hypothesis can be potentially true for other forms of data, more so due to the greater dimensionality of data like image, text, etc. In order to extend the idea to other forms of data, the key will be to find how best to perturb the input to get smooth predictions.

Lastly, we test our hypothesis with LIME and found promising results. Since the instability of explanations of other interpretability methods can also be (at least partly) explained by unstable predictions of the machine learning model, we suspect our idea can be applied to improve other model interpretability methods too.

As we can see, there is a lot of scope for future work and we are excited to see how research develops in this direction.

7. Conclusion

In this paper, we propose a way to improve the robustness of LIME, a model-agnostic explainer of the predictions of machine learning models. We propose smoothening the predictions made by the model to increase the consistency of the predictions made by the model, thereby making the explanations more trustable. We explain how we smoothen predictions using random noise and perform some preliminary experiments on publicly-available datasets to achieve promising results. We also outline future steps that can be taken to increase the scope of the research.

Acknowledgments

We would like to thank Mercari Inc. for supporting the research and also the anonymous reviewers who gave very helpful feedback to improve the quality of the paper.

Any remaining deficiencies left in the paper belong to the authors.

References

- [1] P. Karatza, K. Dalakleidi, M. Athanasiou, K. Nikita, Interpretability methods of machine learning algorithms with applications in breast cancer diagnosis, in: 2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC), 2021, pp. 2310–2313. doi:10.1109/EMBC46164.2021.9630556, ISSN: 2694-0604.
- [2] B. M. Greenwell, B. C. Boehmke, A. J. McCarthy, A Simple and Effective Model-Based Variable Importance Measure, 2018. URL: <http://arxiv.org/abs/1805.04755>. doi:10.48550/arXiv.1805.04755, arXiv:1805.04755 [cs, stat].
- [3] A. Goldstein, A. Kapelner, J. Bleich, E. Pitkin, Peeking Inside the Black Box: Visualizing Statistical Learning with Plots of Individual Conditional Expectation, 2014. URL: <http://arxiv.org/abs/1309.6392>. doi:10.48550/arXiv.1309.6392, arXiv:1309.6392 [stat].
- [4] S. Lundberg, S.-I. Lee, A Unified Approach to Interpreting Model Predictions, 2017. URL: <http://arxiv.org/abs/1705.07874>. doi:10.48550/arXiv.1705.07874, arXiv:1705.07874 [cs, stat].
- [5] S. M. Lundberg, G. G. Erion, S.-I. Lee, Consistent Individualized Feature Attribution for Tree Ensembles, 2019. URL: <http://arxiv.org/abs/1802.03888>. doi:10.48550/arXiv.1802.03888, arXiv:1802.03888 [cs, stat].
- [6] D. Alvarez-Melis, T. S. Jaakkola, On the Robustness of Interpretability Methods, 2018. URL: <http://arxiv.org/abs/1806.08049>. doi:10.48550/arXiv.1806.08049, arXiv:1806.08049 [cs, stat].
- [7] M. T. Ribeiro, S. Singh, C. Guestrin, "Why Should I Trust You?": Explaining the Predictions of Any Classifier (2016). URL: <https://arxiv.org/abs/1602.04938v3>. doi:10.48550/arXiv.1602.04938.
- [8] I. Covert, S.-I. Lee, Improving KernelSHAP: Practical Shapley Value Estimation Using Linear Regression, in: Proceedings of The 24th International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 3457–3465. URL: <https://proceedings.mlr.press/v130/covert21a.html>, ISSN: 2640-3498.
- [9] Y. Zhang, K. Song, Y. Sun, S. Tan, M. Udell, "Why Should You Trust My Explanation?" Understanding Uncertainty in LIME Explanations, 2019. URL: <http://arxiv.org/abs/1904.12991>. doi:10.48550/arXiv.1904.12991, arXiv:1904.12991 [cs, stat].
- [10] Z. Zhou, G. Hooker, F. Wang, S-LIME: Stabilized-LIME for Model Explanation, in: Proceedings

- of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD '21, Association for Computing Machinery, New York, NY, USA, 2021, pp. 2429–2438. URL: <https://doi.org/10.1145/3447548.3467274>. doi:10.1145/3447548.3467274.
- [11] M. R. Zafar, N. Khan, Deterministic Local Interpretable Model-Agnostic Explanations for Stable Explainability, *Machine Learning and Knowledge Extraction* 3 (2021) 525–541. URL: <https://www.mdpi.com/2504-4990/3/3/27>. doi:10.3390/make3030027, number: 3 Publisher: Multidisciplinary Digital Publishing Institute.
- [12] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, R. Fergus, Intriguing properties of neural networks, 2014. URL: <http://arxiv.org/abs/1312.6199>. doi:10.48550/arXiv.1312.6199, arXiv:1312.6199 [cs].
- [13] G. Visani, E. Bagli, F. Chesani, A. Poluzzi, D. Capuzzo, Statistical stability indices for LIME: obtaining reliable explanations for Machine Learning models, *Journal of the Operational Research Society* 73 (2022) 91–101. URL: <http://arxiv.org/abs/2001.11757>. doi:10.1080/01605682.2020.1865846, arXiv:2001.11757 [cs, stat].
- [14] D. Garreau, U. von Luxburg, Looking Deeper into Tabular LIME, 2022. URL: <http://arxiv.org/abs/2008.11092>, arXiv:2008.11092 [cs, stat].
- [15] S. Saito, E. Chua, N. Capel, R. Hu, Improving LIME Robustness with Smarter Locality Sampling, 2021. URL: <http://arxiv.org/abs/2006.12302>. doi:10.48550/arXiv.2006.12302, arXiv:2006.12302 [cs, stat].
- [16] S. Shi, X. Zhang, W. Fan, A Modified Perturbed Sampling Method for Local Interpretable Model-agnostic Explanation, 2020. URL: <http://arxiv.org/abs/2002.07434>. doi:10.48550/arXiv.2002.07434, arXiv:2002.07434 [cs, stat].