

---

# Instructions for Paper Submissions to AISTATS 2023

---

Anonymous Author  
Anonymous Institution

## Abstract

The Abstract paragraph should be indented 0.25 inch (1.5 picas) on both left and right-hand margins. Use 10 point type, with a vertical spacing of 11 points. The **Abstract** heading must be centered, bold, and in point size 12. Two line spaces precede the Abstract. The Abstract must be limited to one paragraph.

## 1 INTRODUCTION

Recently, ML has flourished in critical domains, such as healthcare and finance. In these areas, we need ML models that predict accurately but also with the ability to explain their predictions. Therefore, Explainable AI (XAI) is a rapidly growing field due to the interest in interpreting black box machine learning (ML) models. XAI literature distinguishes between local and global interpretation methods (Molnar et al., 2020). Local methods explain a specific prediction, whereas global methods explain the entire model behavior. Global methods provide a universal explanation, summarizing the numerous local explanations into a single interpretable outcome (number or plot). For example, if a user wants to know which features are significant (feature importance) or whether a particular feature has a positive or negative effect on the output (feature effect), they should opt for a global explainability technique. Aggregating the individual explanations for producing a global one comes at a cost. In cases where feature interactions are strong, the global explanation may obfuscate heterogeneous effects (Herbinger et al., 2022) that exist under the hood, a phenomenon called aggregation bias (Mehrabi et al., 2021).

Feature effect forms a fundamental category of global explainability methods, isolating a single feature’s average impact on the output. Feature effect methods suffer from aggregation bias because the rationale behind the average

effect might be unclear. For example, a feature with zero average effect may indicate that the feature has no effect on the output or, contrarily, it has a highly positive effect in some cases and a highly negative one in others.

There are two widely-used feature effect methods; Partial Dependence Plots (PDPlots)(Friedman, 2001) and Aggregated Local Effects (ALE)(Apley and Zhu, 2020). PDPlots have been criticized for producing erroneous feature effect plots when the input features are correlated due to marginalizing over out-of-distribution synthetic instances. Therefore, ALE has been established as the state-of-the-art feature effect method since it can isolate feature effects in situations where input features are highly correlated.

However, ALE faces two crucial drawbacks. First, it does not provide a way to inform the user about potential heterogeneous effects that are hidden behind the average effect. In contrast, in the case of PDPlots, the heterogeneous effects can be spotted by exploring the Individual Conditional Expectations (ICE)(Goldstein et al., 2015). Second, ALE requires an additional step, where the axis of the feature of interest is split in  $K$  fixed-size non-overlapping intervals, where  $K$  is a hyperparameter provided by the user. This splitting is done blindly, which can lead to inconsistent explanations.

In this paper, we extend ALE with a probabilistic component for measuring the uncertainty of the global explanation. The uncertainty of the global explanation expresses how certain we are that the global (expected) explanation is valid if applied to an instance drawn at random and informs the user about the level of heterogeneous effects hidden behind the expected explanation. Our method completes ALE, as ICE plots complement PDPlots, for revealing the heterogeneous effects.

Our method also automates the step of axis splitting into non-overlapping intervals. We, firstly, transform the bin splitting step into an unsupervised clustering problem and, second, find the optimal bin splitting for a robust estimation of (a) the global (expected) effect and (b) the uncertainty of the explanation from the limited samples of the training set. We formally prove that the objective of the clustering problem has as lower-bound the aggregated uncertainty of the global explanation. Our method works out of the box

without requiring any input from the user.

**Contributions.** The contributions of this paper are the following:

- We introduce Uncertainty DALE (UDALE), an extension of DALE that quantifies the uncertainty of the global explanation, i.e. the level of heterogeneous effects hidden behind the global explanation.
- We provide an algorithm that automatically computes the optimal bin splitting for robustly estimating the explanatory quantities, i.e., the global effect and the uncertainty.
- We formally prove that our method finds the optimal grouping of samples, minimizing the added uncertainty over the unavoidable heterogeneity that is the lower-bound of the objective.
- We provide empirical evaluation of the method in artificial and real datasets.

The implementation of our method and the code for reproducing all the experiments is provided in the submission and will become publicly available upon acceptance.

## 2 BACKGROUND AND RELATED WORK

### 3 THE U-DALE METHOD

#### 3.1 Uncertainty Quantification

The uncertainty of the global explanation emerges from the natural characteristics of the experiment, i.e., the data generating distribution and the black-box function.

##### 3.1.1 Methodology

#### 3.2 Bin Splitting as a Clustering Problem

##### 3.2.1 Methodology

Furthermore, we automate the step of splitting the axis into non-overlapping intervals. The need for non-overlapping bins emerges from the ignorance of the data-generating distribution, enforcing all estimations to be based on the limited instances of the training set. Therefore, there is an implicit trade-off behind the formation of bins. Each bin must include enough instances for a robust estimation of the bin feature effect (expected value), and the uncertainty of the explanation (variance). On the other hand, each bin should include points with similar local effects. Therefore, we transform the bin splitting step into an unsupervised clustering problem, encoding the trade-off mentioned above in

the objective function. We formally show that the objective of the clustering problem has lower-bound the (unavoidable) heterogeneity, i.e., the aggregated uncertainty of the global explanation. Therefore, we aim to find the optimal grouping of samples that adds the slightest uncertainty over the unavoidable heterogeneity. We finally solve the minimization problem by finding the global optimum using dynamic programming. Our method works out of the box without requiring any input by the user. We provide a theoretical and empirical evaluation of our method.

#### 3.2.2 Algorithms

## 4 SYNTHETIC EXAMPLES

## 5 REAL-WORLD EXAMPLES

### Acknowledgements

All acknowledgments go at the end of the paper, including thanks to reviewers who gave useful comments, to colleagues who contributed to the ideas, and to funding agencies and corporate sponsors that provided financial support. To preserve the anonymity, please include acknowledgments *only* in the camera-ready papers.

### References

#### References

- Daniel W Apley and Jingyu Zhu. Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(4):1059–1086, 2020.
- Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil 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.
- Julia Herbinger, Bernd Bischl, and Giuseppe Casalicchio. Repid: Regional effect plots with implicit interaction detection. In *International Conference on Artificial Intelligence and Statistics*, pages 10209–10233. PMLR, 2022.
- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
- Christoph Molnar, Giuseppe Casalicchio, and Bernd Bischl. Interpretable machine learning—a brief history, state-of-the-art and challenges. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 417–431. Springer, 2020.