



Shaping Up SHAP : Enhancing Stability through Layer-Wise Neighbor Selection

Gwladys Kelodjou¹, Laurence Rozé², Véronique Masson¹,
Luis Galárraga¹, Romaric Gaudel¹, Maurice Tchuente³,
Alexandre Termier¹

¹Univ Rennes, Inria, CNRS, IRISA - UMR 6074

²Univ Rennes, INSA Rennes, CNRS, Inria, IRISA - UMR 6074

³Sorbonne University, IRD, University of Yaoundé I, UMI 209 UMMISCO



- Machine Learning techniques are widely used in various domains.
- Concerns about fairness and trustworthiness motivate interest in **explainability**.
 - Addressing machine learning model opacity through **explanations**.

Local post-hoc explainability

Post-hoc local explainability focuses on explaining decisions for single instances.

Local post-hoc explainability approaches :

- LIME¹, SHAP², DeepLIFT³, etc. :
- Providing feature attribution scores.

Some Qualities of a Good Explanation

- **Interpretability** : The degree of understanding of the explanation by a user.
- **Stability** : The ability to reproduce the same explanation.
- **Fidelity** : The ability of the explanation model to accurately mimic the black-box model.

¹Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. «“Why should i trust you?” Explaining the predictions of any classifier». In Proceedings of the 22nd ACM SIGKDD

²Lundberg, S. M.; and Lee, S.-I. 2017. «A unified approach to interpreting model predictions». NIPS 2017

³Shrikumar, A.; Greenside, P.; and Kundaje, A. 2017. «Learning important features through propagating activation differences». In International conference on machine learning, PMLR

- **SHAP (SHapley Additive exPlanations)** : Applies cooperative game theory to reveal feature contributions.
- **Cooperative game theory** : fairly distribute the gain obtained by multiple players collaborating in a game.
- **Shapley value** : Consider different coalitions in which the player is present or absent to determine his contribution.

- **SHAP values** are the **Shapley values** of a conditional expectation function of the original model.
- Exact computation of SHAP values is challenging.
- **Kernel SHAP** is a linear-regression-based approximation of SHAP values.
 - Widely used for its model-agnostic nature.

Problem

Kernel SHAP's estimator **suffers from stability issues.**

- Metrics for stability in LIME explanations are introduced by Visani et al. (2022)¹.
- OptiLIME framework (2020)² allows users to customize the level of stability and fidelity in LIME explanations.
- S-lime (2021)³ proposes a framework to determine the number of perturbation points required to ensure stable LIME explanations.
- The concept of **robustness** is introduced by Alvarez et al. (2018)⁴, highlighting challenges in popular methods like LIME and SHAP.

¹Visani, G. et al. 2022. «Statistical stability indices for LIME : Obtaining reliable explanations for machine learning models». Journal of the Operational Research Society

²Visani, G.; Bagli, E.; and Chesani, F. 2020. **OptiLIME : Optimized LIME Explanations for Diagnostic Computer Algorithms**. CIKM Workshops

³Zhou, Z.; Hooker, G.; and Wang, F. 2021. «S-lime : Stabilized-lime for model explanation». KDD

⁴Alvarez-Melis, D.; and Jaakkola, T. S. 2018. «On the Robustness of Interpretability Methods». ICML Workshop

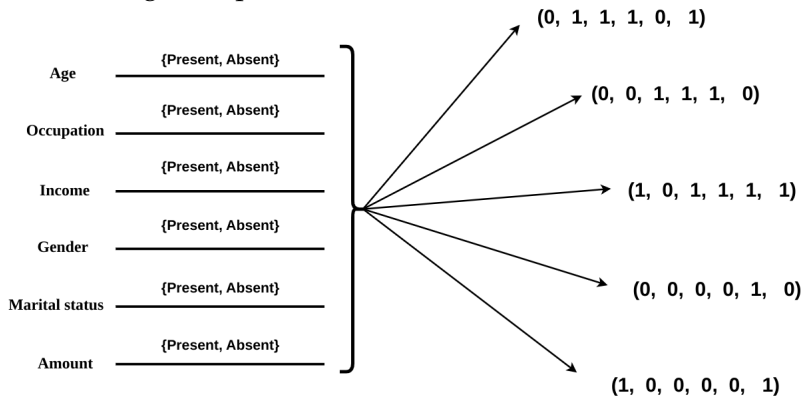
Coalition : subset of features.

Age	Present	1
Occupation	Present	1
Income	Absent	0
Gender	Present	1
Marital status	Absent	0
Amount	Present	1

(1, 1, 0, 1, 0, 1)

A coalition is also referred to as a **neighbor**.

Determining the contributions of each feature requires considering multiple coalitions.

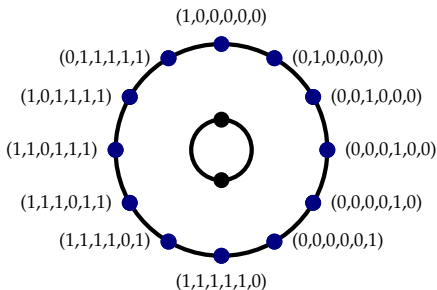


Computing SHAP values involves considering all possible coalitions.

Layer-Wise Neighbor Selection

Layer

Set of coalitions sharing the same number of features present or features absent.



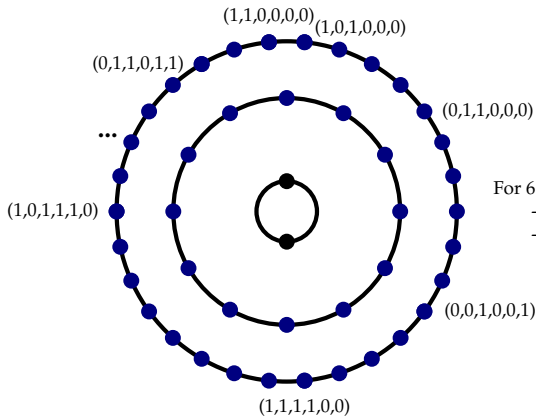
For 6 features, **layer 1** contains **12 coalitions** :

- All coalitions with a single feature present
- All coalitions with a single feature absent

Layer 1

Coalitions with a single feature present or with a single feature absent.

Layer-Wise Neighbor Selection

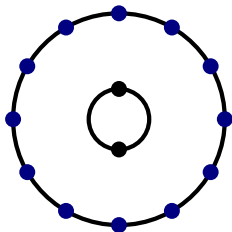


For 6 features, **stratum 2** contains **42 coalitions** :

- Coalitions of layer 1 (12)
- Coalitions of layer 2 (30)

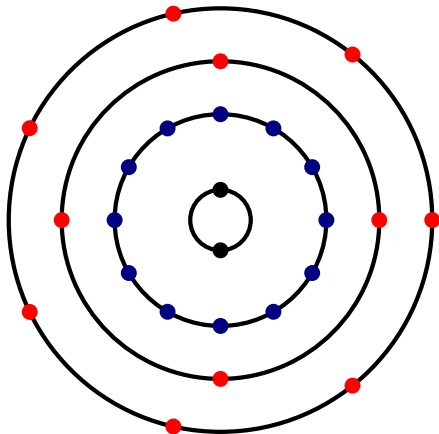
- **Layer 2** : Coalitions with a two features present or two features absent.
- **Stratum 2** : Layer 1 + Layer 2.

- Samples a subset of coalitions.
- **Budget** : determines the size of the coalition sample.



Stratum 1 full !

First generate coalitions from lower layers.

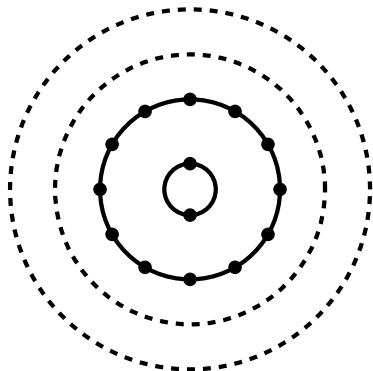


Random sampling within layers 2 and 3

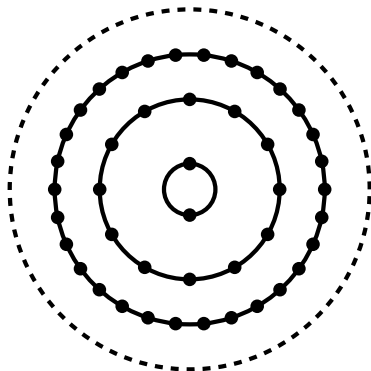
Randomly samples from subsequent layers if the budget is not exhausted.

Achieving Kernel SHAP's stability

Set the budget to consider only complete stratum.

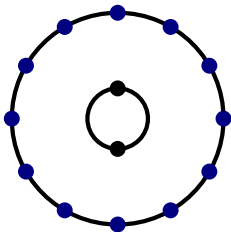


Stratum 1 full



Stratum 2 full

Contribution 2 : Stratum 1 Attributions



Use only Stratum 1 coalitions

Attribution values with Stratum 1

For any feature $j \in N = \{1, \dots, M\}$ (the set of all features), attribution ϕ_j with Stratum 1 is :

$$\phi_j = \tilde{\phi}_j + \frac{1}{M} \left(f(N) - f(\emptyset) - \sum_{i=1}^M \tilde{\phi}_i \right)$$

where for any i , $\tilde{\phi}_i = \frac{f(\{i\}) - f(\emptyset) + f(N) - f(N \setminus \{i\})}{2}$ and M the number of features.

Stratum 1 Attribution Properties

- **LES family^{2,3} :**
 - Linearity
 - Efficiency
 - Symmetry
- Missingness
- **Execution time : $O(M)$**

SHAP values Properties

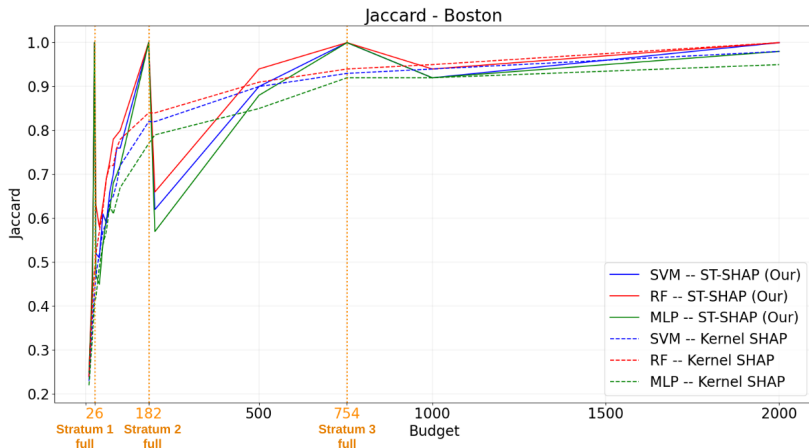
- Local Accuracy
 - Efficiency
- Missingness
- Consistency
 - Null effect
 - Linearity
 - Symmetry
- **Execution time : $O(2^M)$ ¹**

¹ M is the number of features in the example to be explained.

² Ruiz, L. M.; Valenciano, F.; and Zarzuelo, J. M. 1998. «The family of least square values for transferable utility games». Games and Economic Behavior

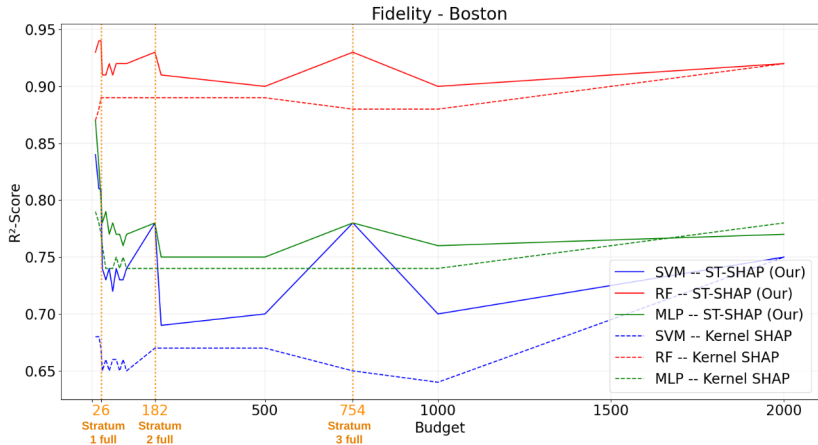
³ Condevaux, C.; Harispe, S.; and Mussard, S. 2022. Fair and Efficient Alternatives to Shapley-based Attribution Methods. ECML PKDD

- **Black-box models** : Support Vector Machine (SVM), Random Forest (RF), Logistic Regression (LR), Multi-layer Perceptron (MLP).
- **Datasets** : Boston, Movie, Credit Card, Adult Income, Dry Bean, HELOC, Spambase, Wisconsin Breast Cancer Database.
- **Various budgets** to assess stability and fidelity.
- **Metrics** :
 - **Jaccard coefficient** for stability measurement.
 - R^2 -**score** and **Accuracy** for fidelity measurement.
- Code is publicly available at <https://github.com/gwladyskelodjou/st-shap>.



Complete stratums \Rightarrow stable output.

Experiments Results



No impact on explanation fidelity observed with our methods :
remains highly accurate.

- Eliminating the random step in Kernel SHAP leads to explanation stability.
- Removing randomness still maintains high fidelity.
- Using only Stratum 1 achieves complete stability and good fidelity.

- Properties of black-box models that guarantee good approximations of the SHAP values.
- Budget - Black-Box complexity Relationship for Explanation Optimization.