

Shapley Values based Regional Feature Importance Measures driving Error Analysis in Manufacturing

Introduction

HOCHSCHULE FURTWANGEN UNIVERSITY



Introduction



Feature importance measures driving error analysis in manufacturing

- Data driven error analysis can leverage predictive models for error cause analysis
 - Quality Management analyzes features pointing at "interesting" phenomena
 - Feature importance measures provide insights to the Human-In-the-Loop
 - State-of-the-Art feature importance measures are not tailored to this task

Contribution:

Furtwangen University

- Define an applied notion of "interestingness" and proposed three approaches to determine this as feature importance measure
 - Evaluating and comparing with state-of-the-art importance measures on synthetic and real-world data



3

Agenda



- Motivation and Use Case
- Challenges with feature importance in error analysis
- Proposed SHAP-based approaches to determine regional feature importance
- Experiments and results
- Conclusion

Motivation and Use Case

Explaining errors in manufacturing data driven quality processes

- Proactive error prevention for rare but costly errors in production [1]
- Data driven error analysis approach to determine error causes
- Domain experts evaluate production data and take corrective steps in production
- Traditionally, cause analysis is performed by evaluating input features with correlated errors
- Increasing number of features requires
 automated feature ranking methods
- Which measures of feature importance suits the task for cause analysis?



Histograms of features with highlighted errors traditionally used for cause analysis [1]. Color coded error percentage: white=0%, black=100%



Challenges with feature importance



Global importance



Furtwangen University

Feature Importance

Challenges with feature importance in cause analysis



Holistic view:

- Importance based on global "averages"
- Missing local phenomena
- Examples: Gain, Weight/Frequency, Gini, Abs. Avg. SHAP, ...



Global

Local

Regional

Feature Importance

Challenges with feature importance in cause analysis



Global **Local** Regional

Instance based:

- Importance based on single samples
- Isolated samples provide little general insights
- Missing context e.g., relations to other samples



Challenges with feature importance in cause analysis



What is needed for feature importance in cause analysis?

A feature should have high importance if it **contributes to interesting predictions** [1]:

- *it provide at least sometimes strong hints for errors*
- *it is relevant in at least some cases*
- it allows to draw conclusions upon inspection



Global and generalizing



Holistic but targeted to interesting aspects



Case specific

Global Local

Regional



Shapley values as basis for regional feature importance



Shapley values as basis for importance measures

Feature importance for cause analysis

"The Shapley value is the average marginal contribution of a feature value across all possible coalitions." [3]

SHAPley Values

Tree Explainer

Global Importance

- Desirable properties: Accuracy, consistency, missingness
- Shapley values can be computed for each prediction individually
- SHAP (SHapley Additive exPlanation)





Feature importance for cause analysis

 $\in \mathbb{R}$ of feature f on model m is attributed using Shapley **SHAPley Values** $\phi_f = \sum \frac{|S|! (M - |S| - 1)!}{M!} [m_x(S \cup \{f\}) - m_x(S)]$ **Global Importance**

where *M* is the number of all features and S the set of input values.



The contribution
$$\phi_f$$
 evalues [3]:

Shapley values as basis for importance measures

Feature importance for cause analysis

The contribution $\phi_f \in \mathbb{R}$ of feature f on model m is attributed using Shapley values [3]:

$$\phi_f = \sum_{S \subseteq M \setminus \{f\}} \frac{|S|! (M - |S| - 1)!}{M!} [m_x(S \cup \{f\}) - m_x(S)]$$

SHAPley Values

where M is the number of all features and S the set of input values.

Tree Explainer

Global Importance



Figure – "SHAP feature attribution ": Contribution of each feature as change in the expected model prediction when conditioning on that feature (source [4]).



Shapley values as basis for importance measures

Feature importance for cause analysis

- Implementation for trees [4] (i.e., XGBoost)
- Computation of exact Shapley values in polynomial time
- O(TLD²), where T is the number of trees, L is the maximum number of leaves in any tree and D the maximal depth of any tree

SHAPley Values Tree Explainer

Global Importance





Feature importance for cause analysis

Idea: Features with large absolute Shapley values are important

SHAPley Values

Tree Explainer Global Importance Average absolute Shapley values per feature across the data [4]:

$$I_f = \frac{1}{n} \sum_{i=1}^{n} |\phi_f^{(i)}|$$

HOCHSCHULE FURTWANGEN UNIVERSITY



Importance measures that reveal insights on errors

Previous Work

Importance measures that reveal insights on errors [1]



Fundamental idea:

• Aggregate Shapley values so that "interesting" features get a high score (i.e., not just averaging them for a global)

Name	Intuition	Formal Definition
Max Shap	Highest Shapley value across analysed dataset	$Max SHAP_f(m, S) = \max\{\phi_f(m, x) x \in S\}$
Max Main	Like Max SHAP but without interaction effects	$Max \ Main \ Effect_f(m,S) = \max\{\phi_f(m,x) - \sum_{j \neq f} \phi_{f,j}(m,x) \ x \in S\}$
Range Shap	Range of Shapley values across analysed dataset	Range SHAP _f (m, S) = max{ $\phi_f(m, x) x \in S$ } - min{ $\phi_f(m, x) x \in S$ }



Proposed SHAP-based regional feature importance measures

Proposed Regional Feature Importance Measures



Outlier-Approach

Micro-Average Approach Slope-Approach

Fundamental idea:

• A scoring-function $g: g(f, X, ...) \rightarrow \mathbb{R}$ that aggregates SHAP values and scores "interesting" features high

f : target feature

X: dataset

Proposed regional feature importance measures



Idea: A feature with abnormal Shapley values are potentially interesting

Outlier-Approach

Micro-Average Approach Perform anomaly detection over the distribution of SHAP values: Slope-Approach

$$g(f, X, \lambda) = \sum_{x' \in \text{outl}(\lambda, X)} \phi_f(x')$$

f: target featureX: dataset σ : standard deviation λ : multiplier of σ $\overline{\phi}_f$: mean shap value

outl(
$$\lambda, X$$
) = { $x \in X | \phi_f(x) \ge \overline{\phi}_f(X) + \lambda \sigma(\phi_f(X))$ }

Proposed regional feature importance measures



Idea: A feature is of interest if it shows high Shapley values within a small feature value range

Outlier-Approach

Micro-Average Approach Slope-Approach

Partition the feature and determine the average SHAP values over equally sized intervals:

$$g(f, X, n) = \max\{\overline{\phi}_f(X_i) \mid i = 0, \dots, n-1\}$$

f: target feature *X*: dataset *n*: number of intervals $\overline{\phi}_f$: mean shap value *X_i*: feature interval

$$X_i = \{x \in X | (i * d) \le x < (i * d) + d\}$$

Proposed regional feature importance measures



Idea: A feature with rapid changes in Shapley values is of interest

Outlier-Approach

Micro-Average Approach Slope-Approach

Partition the feature and determine regression slopes over the means of SHAP values over the intervals:

$$g(f, X, n, w) = max\{|slope(w_j)|| j = 0, ..., n - 1\}$$

f: target feature X: dataset n: number of intervals w: window size $\overline{\phi}_f$: mean shap value X_i : feature interval

$$w_{j} = \{ \bar{\phi}_{f}(X_{j}), \dots, \bar{\phi}_{f}(X_{j+w}) \}$$
$$\bar{\phi}_{f}(X_{i}) = i * \beta + \epsilon \mid \forall \bar{\phi}_{f}(X_{i}) \in w_{j}$$



Experiments

Furtwangen University



Setup 1: "Synthetic Dataset"

Experiments with synthetic data

- Creation of simple data sets with known "*interestingness*" of features
- Comparison of proposed importance measures with established importance measures against know ground truth

Setup 2: "Real-World Dataset"

Experiments with real-world manufacturing data

- Scoring of feature with proposed and established importance measures
- Manual inspection of features regarding *"interestingness"*

- Binary classification with XGBoost with training score f1=1.0
- Computation of importance score on the training set

Experiments synthetic data

Example experiment with synthetic data

- Features A, B with e.g., uniform distribution 0 to 1 and one feature noise
- 3% of data points have errors
- If 0.4 < B < 0.425, then 8% of data points have errors
- If A < 0.6, then 2.5% of data points have errors





HOCHSCHULE FURTWANGEN UNIVERSITY

Figure "Shap plots of feature A and B": Not interesting feature A (left) and interesting feature B (right). Color coded – red: error, blue: non error, green marker: area of interest

|--|

Interesting

R	Classic Metrics						SHAP-based Metrics				Proposed Metrics			
а						Average			Range	Slope App	roaches	Micro-Ava	Outlier	
n k	Weight	Gain	Cover	Total Gain	Total Cover	Abs SHAP	Max Main	Max Shap	SHAP	Max SHAP	IQR SHAP	Approach	Approach	
1	A	A	A	A	A	A	A	A	A	В	В	В	В	
2	В	В	В	В	В	В	noise	noise	В	Noise	A	A	A	
3	noise	noise	noise	noise	noise	noise	В	В	noise	A	noise	noise	noise	

Experiment with real world-data

- Secom dataset originating from a semiconductor manufacturer [6]
- 591 features and 1667 instances with 106 error instances
- Manual assessment of the top five features of all importance measure
- Comparison of grouped metrics with proposed metrics

R	Classic Metrics					SHAP-based Metrics				Proposed Metrics			
a n k	Weight	Gain	Cover	Total Gain	Total Cover	Average Abs Shap	Max Main	Max Shap	Range Value	Slope App Max Shap	oroaches IQR Shap	Micro-Avg Approach	Outlier Approach
1	F59	F210	F168	F59	F59	F59	F59	F59	F59	F59	F59	F64	F59
2	F333	F539	F429	F333	F64	F21	F64	F64	F64	F64	F64	F59	F423
3	F103	F29	F426	F64	F426	F333	F40	F426	F333	F429	F33	F103	F64
4	F2	F109	F100	F132	F121	F488	F426	F333	F103	F333	F130	F40	F333
5	F33	F304	F331	F33	F574	F103	F153	F40	F33	F475	F429	F121	F2



Experiment with real world-data

0.0

-0.2

50





3

4 F475 -0.8

-1.0

0.5

0.6

0.7

F130

0.8

0.9

0.0

200

150

100

F423

Agreements of importance

Disagreements of importance

Highlights of proposed

Furtwangen University



Conclusion:

- Shapley values have potential as the basis for regional importance measures
- Regional feature importance measures can pinpoint interesting features for cause analysis in manufacturing
- Experiments show the usefulness of proposed SHAP-based approaches for cause analysis

Future Work:

- Evaluation on "interestingness" needs more objective measure
- Evaluation of importance measures with domain experts





REFERENCES



[1] H. Ziekow, U. Schreier, A. Gerling, and A. Saleh, "Interpretable Machine Learning for Quality Engineering in Manufacturing - Importance Measures that Reveal Insights on Errors", *The Upper-Rhine Artificial Intelligence Symposium, UR-AI 2021, Artificial Intelligence - Application in Life Sciences and Beyond*, Germany, Kaiserslautern: Hochschule Kaiserslautern, University of Applied Sciences, pp. 96–105, October 2021

[2] XGBoost Documentation, "Python API", *Reference. xgboost developers*. [Online]. Available from: https://xgboost.readthedocs.io/en/latest/python/python_api.html, retrieved on 08/26/2022

[3] C. Molnar, "Interpretable Machine Learning: A Guide for Making Black Box Models Explainable", *2nd edn.*, 2022. [Online]. Available from: https://christophm.github.io/interpretable-ml-book, retrieved on 08/26/2022.

[4] S. M. Lundberg, and S. I. Lee, "A unified approach to interpreting model predictions", *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17)*, NY, USA, pp. 4768–4777, 2017.

[5] S. M. Lundberg et al., "From local explanations to global understanding with explainable AI for trees", *Nature machine intelligence*, *2(1)*, pp. 56–67, 2020.

[6] D. Dua and C. Graff, "UCI Machine Learning Repository", Irvine, CA: University of California, School of Information and Computer Science, 2019. [Online]. Available from: http://archive.ics.uci.edu/ml.