



# Explainable Artificial Intelligence

## Modell-agnostische Erklärungsansätze im Vergleich

Büsra Karaoglan

Fachbereich Mathematik, Naturwissenschaften und Datenverarbeitung  
Studiengang Business Mathematics

29. September 2020

# Inhaltsverzeichnis

- 1 Motivation
- 2 Grundlagen der XAI
- 3 Local Interpretable Model-Agnostic Explanations
- 4 Shapley Additive Explanations
- 5 Implementierung und Evaluation

# Inhaltsverzeichnis

- 1 Motivation
- 2 Grundlagen der XAI
- 3 Local Interpretable Model-Agnostic Explanations
- 4 Shapley Additive Explanations
- 5 Implementierung und Evaluation

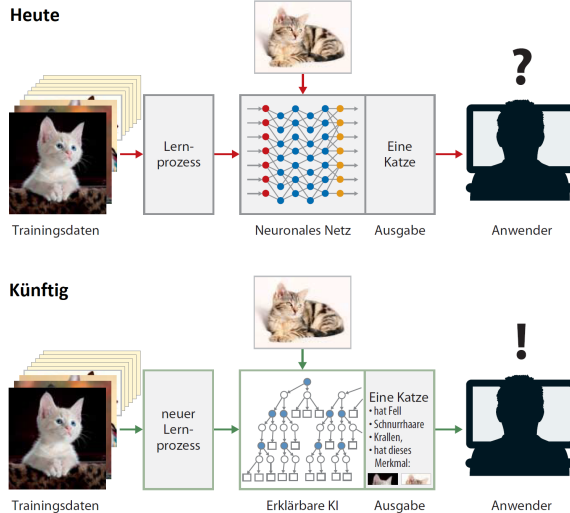


Abbildung: Konzept der erklärbaren Künstlichen Intelligenz (HA18)

# Inhaltsverzeichnis

- 1 Motivation
- 2 Grundlagen der XAI
- 3 Local Interpretable Model-Agnostic Explanations
- 4 Shapley Additive Explanations
- 5 Implementierung und Evaluation

- Clever-Hans-Effekt
- Scheinkausalität
- Algorithmischer Voreingenommenheit

# Inhaltsverzeichnis

- 1 Motivation
- 2 Grundlagen der XAI
- 3 Local Interpretable Model-Agnostic Explanations
- 4 Shapley Additive Explanations
- 5 Implementierung und Evaluation

$$\xi(x) = \operatorname{argmin}_{g \in G} L(f, g, \pi_x(z)) + \Omega(g)$$

- $d$ -dimensionaler Merkmalsraum  $X = \mathbb{R}^d$
- Ausgaberaum  $Y = \mathbb{R}$
- Black-Box-Modell  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  mit  $x \in \mathbb{R}^d$
- Interpretierbares Modell  $g : \mathbb{R}^{d'} \rightarrow \mathbb{R}$
- Komplexität  $\Omega(g)$
- Transformationsfunktion  $IR : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  mit  $x' \in \mathbb{R}^{d'}$
- Ähnlichkeitsmaßes  $\pi_x(z)$  mit permutierte Instanzen  $z' \in \mathbb{R}^{d'}$
- Verlustfunktionen  $L(f, g, \pi_x(z))$



---

## Algorithm 1 Spärliche lineare Erklärungen mit LIME

---

**Require:** Black-Box Modell  $f$ , Stichprobengröße  $N$

**Require:** Instanz  $x$ , dazugehörige interpretierbare Darstellung  $x'$

**Require:** Ähnlichkeitsmaß  $\pi_x(z)$ , Länge der Erklärung  $K$

```
1:  $\mathcal{Z} \leftarrow \{\}$ 
2: for  $i \in \{1, 2, 3, \dots, N\}$  do
3:    $z'_i \leftarrow \text{sample\_around}(x')$ 
4:    $\mathcal{Z} \leftarrow \mathcal{Z} \cup \langle z'_i, f(z_i), \pi_x(z_i) \rangle$ 
5: end for
6:  $w \leftarrow \text{K-LASSO}(\mathcal{Z}, K)$     ▷ mit  $z'_i$  als Merkmale,  $f(z)$  als Zielvariable
7: return  $w$ 
```

---

# LIME-Algorithmus für Tabellendaten

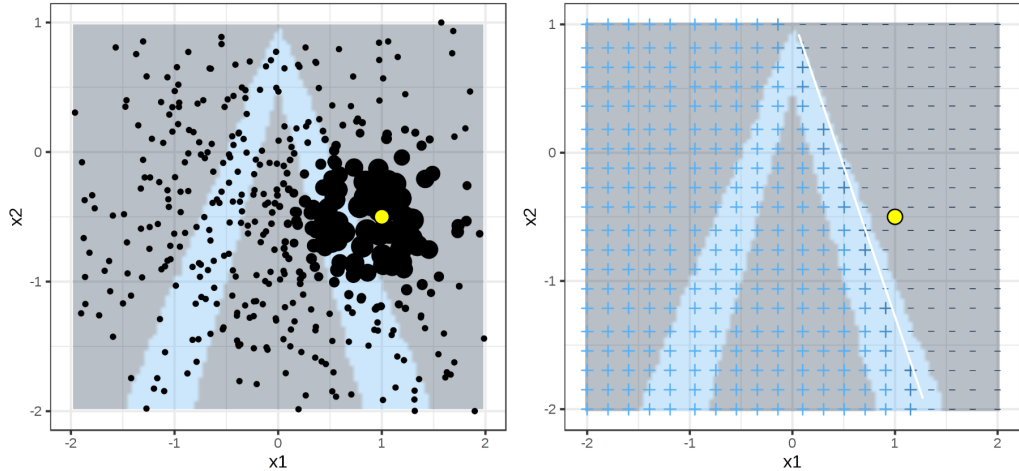


Abbildung: LIME-Algorithmus für tabellarische Daten (MC20)

# Inhaltsverzeichnis

- 1 Motivation
- 2 Grundlagen der XAI
- 3 Local Interpretable Model-Agnostic Explanations
- 4 Shapley Additive Explanations**
- 5 Implementierung und Evaluation

Shapley-Wert für kooperative Spiele:

$$\phi_j(v) = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|! (N - |S| - 1)!}{N!} (v(S \cup \{j\}) - v(S))$$

Shapley-Wert für erklärbares maschinelles Lernen:

$$\phi_j = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|! (M - |S| - 1)!}{M!} [f_x(S \cup \{j\}) - f_x(S)]$$

Additive Merkmalszuordnung:

$$g(z') = \phi_0 + \sum_{j=1}^M \phi_j z'_j$$

SHAP-Wert:

$$\phi_j(f, x) = \sum_{z' \subseteq x'} \frac{|z'|! (M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus j)]$$

$$\Omega(g) = 0$$

$$\pi_{x'}(z') = \frac{(M-1)}{\binom{M}{|z'|} |z'| (M - |z'|)}$$

$$L(f, g, \pi_{x'}) = \sum_{z' \in Z} [f(h_x(z')) - g(z')]^2 \pi_{x'}(z')$$

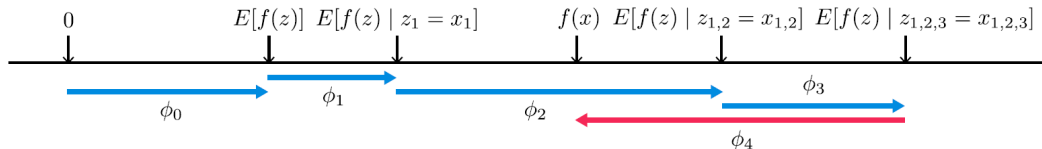


Abbildung: Schematische Darstellung der SHAP-Werte (LL17)

# Inhaltsverzeichnis

- 1 Motivation
- 2 Grundlagen der XAI
- 3 Local Interpretable Model-Agnostic Explanations
- 4 Shapley Additive Explanations
- 5 Implementierung und Evaluation

## Default Of Credit Card Clients Datensatz:

- ID
- LIMIT\_BAL
- SEX
- EDUCATION
- MARRIAGE
- PAY\_0 - PAY\_6
- BILL\_AMT1 - BILL\_AMT6
- PAY\_AMT1 - PAY\_AMT6
- default payment next month



Modell	Precision	Recall-Wert	F1-Wert	Accuracy	ROC
SVM Classifier	0.52	0.52	0.52	0.79	0.69
RF Classifier	0.52	0.51	0.52	0.79	0.69
MLP Classifier	0.42	0.63	0.53	0.73	0.70

**Tabelle:** Übersicht der Klassifizierungsmetriken von angewendeten Modellen

# LIME- und Kernel-SHAP-Erklärungen I

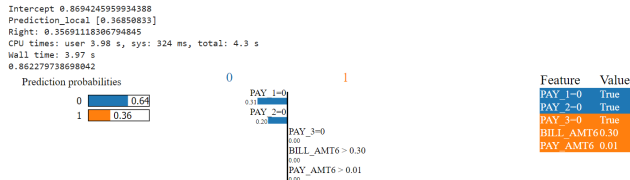


Abbildung: LIME-Erklärung zum SVM Classifier für Instanz 1

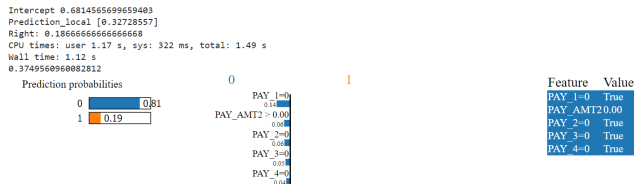
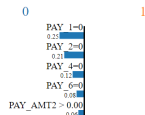


Abbildung: LIME-Erklärung zum RF Classifier für Instanz 1

```
Intercept 1.011229539754286
Prediction_local [0.2923108]
Right: 0.2066158066500271
CPU times: user 979 ms, sys: 250 ms, total: 1.23 s
Wall time: 930 ms
0.44401059872943155
```

Prediction probabilities



Feature	Value
PAY_1=0	True
PAY_2=0	True
PAY_4=0	True
PAY_6=0	True
PAY_AMT2 0.00	

Abbildung: LIME-Erklärung zum MLP Classifier für Instanz 1

```
CPU times: user 43.7 s, sys: 1.74 s, total: 45.4 s
Wall time: 41.6 s
```



Abbildung: Kernel-SHAP-Erklärung zum MLP Classifier für Instanz 1

CPU times: user 1min 22s, sys: 8.07 s, total: 1min 30s  
Wall time: 1min 21s  
[<Figure size 432x288 with 1 Axes>, <Figure size 432x288 with 1 Axes>]

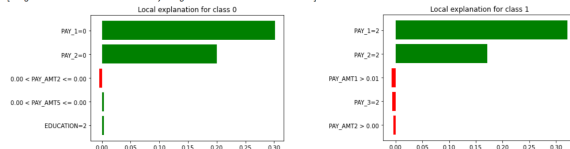


Abbildung: SP-LIME-Erklärung zum SVM Classifier

CPU times: user 24.1 s, sys: 8.19 s, total: 32.3 s  
Wall time: 22.8 s  
[<Figure size 432x288 with 1 Axes>, <Figure size 432x288 with 1 Axes>]

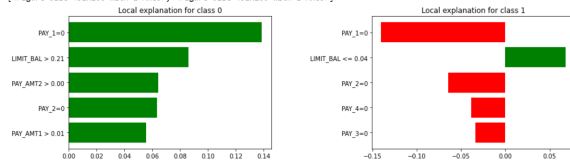


Abbildung: SP-LIME-Erklärung zum RF Classifier

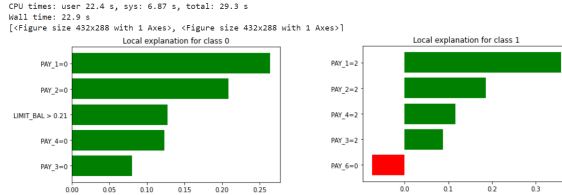


Abbildung: SP-LIME-Erklärung zum MLP Classifier

# LIME- und Kernel-SHAP-Erklärungen V

CPU times: user 30min 48s, sys: 18min 35s, total: 49min 23s  
Wall time: 26min 35s

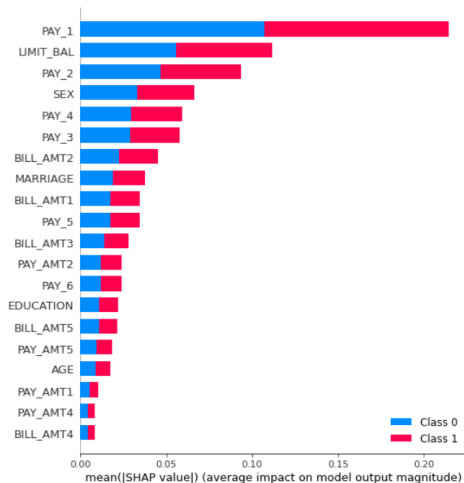


Abbildung: Globale Kernel-SHAP-Erklärung zum MLP Classifier

- Holzinger, A.  
Interpretierbare KI - Neue Methoden zeigen Entscheidungswege künstlicher Intelligenz auf  
Heise Medien, Heft 22, 2018
- Lundberg, S. M. und Lee, S.-I.  
A unified approach to interpreting model predictions  
In: Advances in Neural Information Processing Systems, 2017
- Molnar, C.  
Interpretable Machine Learning - A Guide for Making Black Box Models Explainable  
Abgerufen am 20.08.2020 von  
<https://christophm.github.io/interpretable-ml-book/>
- Ribeiro, M. T., Singh, S. und Guestrin, C.  
„Why Should I Trust You?“ - Explaining the Predictions of Any Classifier  
In: Knowledge Discovery and Data Mining, 2016

Danke für Ihre Aufmerksamkeit