

Can We Trust a Neural Network Prediction?

Methods and Pitfalls for Explaining Black Boxes

Niklas Koenen¹² & Marvin N. Wright¹²³

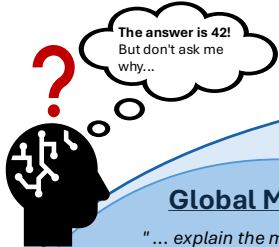
¹Leibniz Institute for Prevention Research & Epidemiology – BIPS

²Faculty of Mathematics and Computer Science, University of Bremen

³Department of Public Health, University of Copenhagen

December 9th, 2024

32nd International Biometric Conference – Atlanta



eXplainable AI (XAI)

Global Methods

"... explain the model's overall behavior across the entire dataset."



Accumulated Local Effects (ALE)
— Apley & Zhu (2020)

Partial Dependence Plots (PDP)
— Friedman (2001)

Permutation Feature Importance (PFI)
— Fisher et al. (2019)

SAGE

— Covert et al. (2020)

Functional ANOVA

— Hooker (2004)

Local Methods

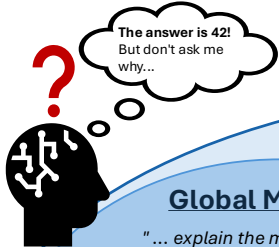
"... explain specific predictions or outcomes for individuals."



Local Surrogate (LIME)
— Ribeiro et al. (2016)

Counterfactual Expl.
— Wachter et al. (2017)

ICE — Goldstein et al. (2015)



eXplainable AI (XAI)

Global Methods

"... explain the model's overall behavior across the entire dataset."



Accumulated Local Effects (ALE)
— Apley & Zhu (2020)

Partial Dependence Plots (PDP)
— Friedman (2001)

Permutation Feature Importance (PFI)
— Fisher et al. (2019)

SAGE
— Covert et al. (2020)

Functional ANOVA
— Hooker (2004)

Local Methods

"... explain specific predictions or outcomes for individuals."



Local Surrogate (LIME)
— Ribeiro et al. (2016)

Counterfactual Expl.
— Wachter et al. (2017)

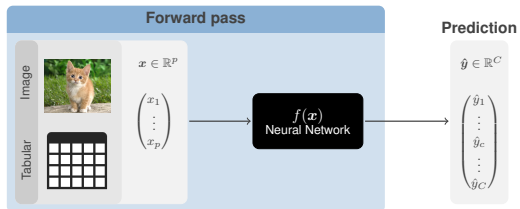
ICE — Goldstein et al. (2015)

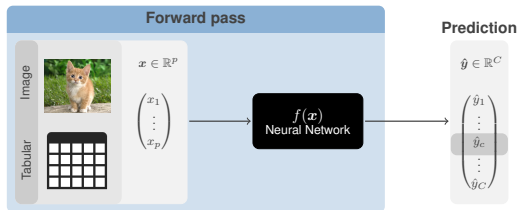
Feature Attribution

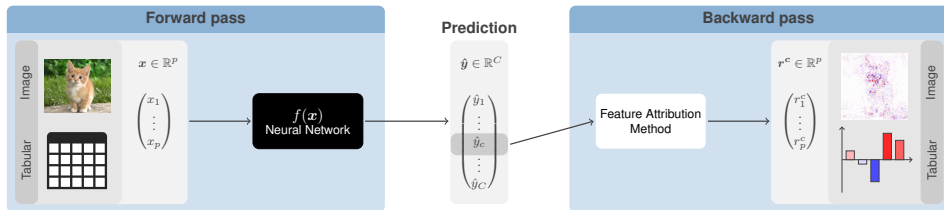
Shapley
SHAP
— Lundberg & Lee (2017)

Neural Networks

From local to global (?)

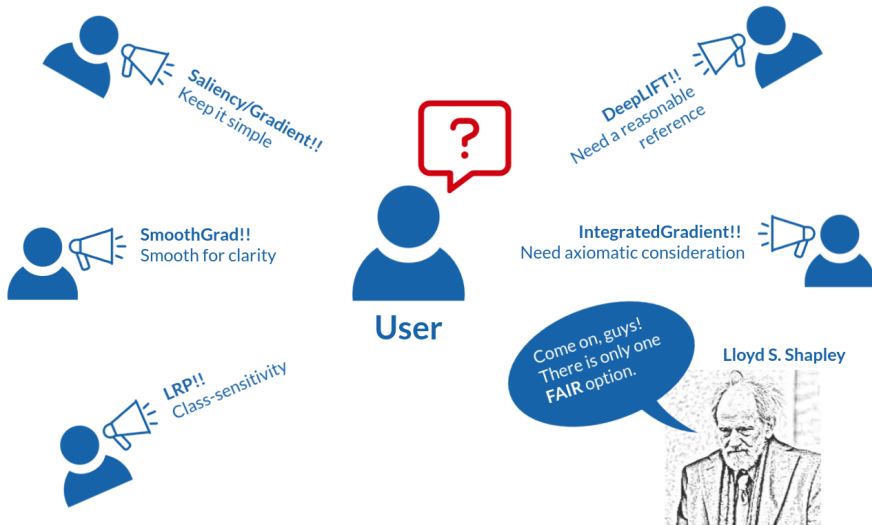






- Assigns an **attribution value** of a selected output to each input variable
 - also known as contribution or relevance
 - can be positive (in **red**) or negative (**blue**)
- Utilizes the
 - layer-wise architecture of a neural network
 - automatic differentiation engine of the training
- Extremely efficient and fast
 - thanks to deep learning libraries like PyTorch, Keras/TensorFlow etc.
 - generally only one forward and one backward pass is needed (no optimization or estimation)

Which Method Should I Use?



Which Method Should I Use?



4

Published in Transactions on Machine Learning Research (06/2024)

The Disagreement Problem in Explainable Machine Learning: A Practitioner's Perspective

Satyapriya Krishna^{1,*}, Tessa Han^{1,*}, Alex Gu², Steven Wu³, Shahin Jabbari⁴, Himabindu Lakkaraju¹

¹*Harvard University*

²*Massachusetts Institute of Technology*

³*Carnegie Mellon University*

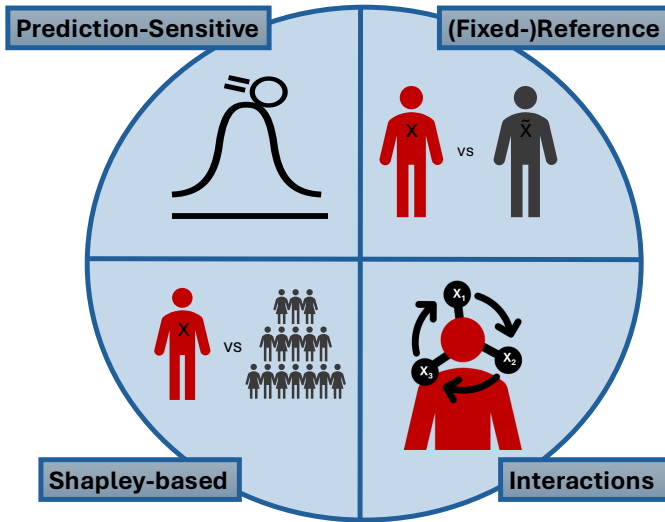
⁴*Drexel University*

* These authors contributed equally to this work.

Reviewed on OpenReview: <https://openreview.net/forum?id=jESY2WTZCe>



Which Method Should I Use? – It depends...



Simulation

Features (real):

BMI, Age, Gender and HbA1c

Outcome (binary)

Diabetes (synth.)



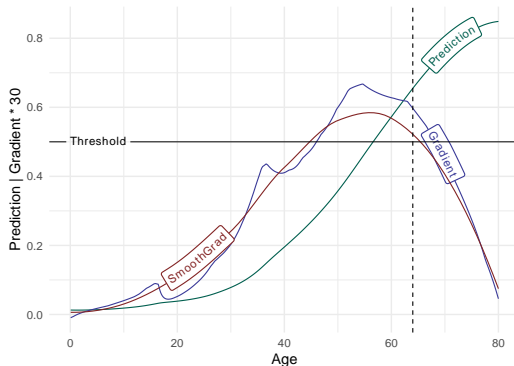
"How sensitive is the prediction
w.r.t. this feature?"

Gradient¹: $\frac{\partial f(x)}{\partial x_i}$ **SmoothGrad²:** $\sum_{k=1}^K \frac{\partial f(x + \varepsilon^{(k)})}{\partial x_i + \varepsilon_i^{(k)}} \quad (\varepsilon^{(k)} \sim \mathcal{N}(0, I\sigma))$

6

Interpretation:

- Increasing patient's age increases model's prediction for diabetes



¹Simonyan et al. (2014) • ²Smilkov et al. (2017)



"How sensitive is the prediction
w.r.t. this feature?"

Gradient¹: $\frac{\partial f(x)}{\partial x_i}$ **SmoothGrad²:** $\sum_{k=1}^K \frac{\partial f(x + \epsilon^{(k)})}{\partial x_i + \epsilon_i^{(k)}} \quad (\epsilon^{(k)} \sim \mathcal{N}(0, I\sigma))$

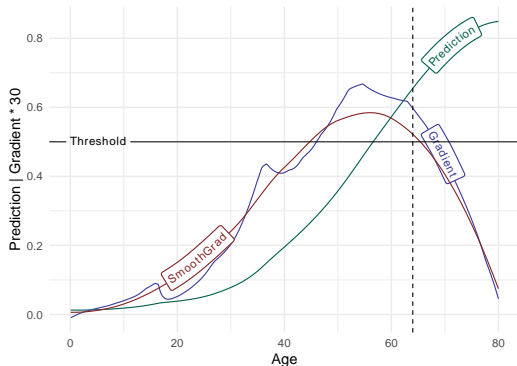
6

Interpretation:

- Increasing patient's age increases model's prediction for diabetes

Pitfalls:

- Show output sensitivity, not attributions!
- Depend on feature scaling
→ complicating comparisons for tabular data
- Provide point-specific insights
- **SmoothGrad**
 - Reduces noise ("area-gradient") but depends on smoothing parameters
 - Is Gaussian noise always the best choice?



¹Simonyan et al. (2014) • ²Smilkov et al. (2017)



"What are the features' contributions for predicting $f(x)$ instead of $f(\tilde{x})$?"

*Grad×Input³: $\frac{\partial f(x)}{\partial x_i} \cdot x_i$ IntGrad⁵: $(x_i - \tilde{x}_i) \int_{\alpha=0}^1 \frac{\partial f(\tilde{x} + \alpha(x - \tilde{x}))}{\partial x_i} d\alpha$

*LRP⁴

DeepLIFT³

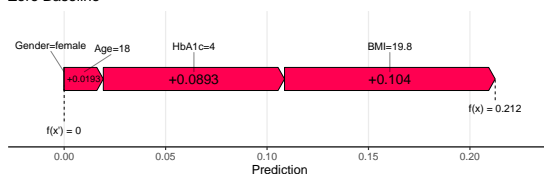
*Approximation and fixed reference ($\tilde{x} = 0$)

7

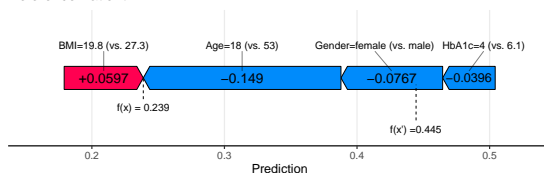
Interpretation:

- Decomposition of $f(x) - f(\tilde{x})$ in feature-wise effects
- Did this feature change (from x to \tilde{x}) argue against or for an increase?

Zero Baseline



Reference Patient





"What are the features' contributions for predicting $f(x)$ instead of $f(\tilde{x})$?"

*Grad×Input³: $\frac{\partial f(x)}{\partial x_i} \cdot x_i$ IntGrad⁵: $(x_i - \tilde{x}_i) \int_{\alpha=0}^1 \frac{\partial f(\tilde{x} + \alpha(x - \tilde{x}))}{\partial x_i} d\alpha$

*LRP⁴

DeepLIFT³

*Approximation and fixed reference ($\tilde{x} = 0$)

7

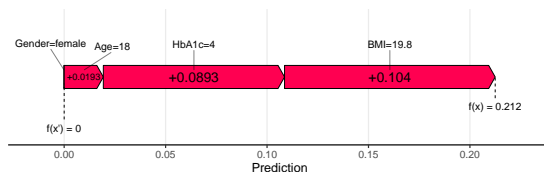
Interpretation:

- Decomposition of $f(x) - f(\tilde{x})$ in feature-wise effects
- Did this feature change (from x to \tilde{x}) argue against or for an increase?

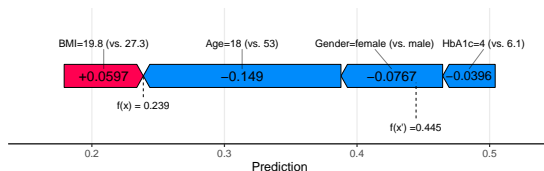
Pitfalls:

- Explanations are highly sensitive to the choice of reference
- Other reference means other question
- Grad×Input and LRP struggle with non-linearity⁶
- Typical references (e.g., 0) often fall outside data distribution

Zero Baseline



Reference Patient





"What are the features' contributions for predicting $f(x)$ compared to $\mathbb{E}_X[f(X)]$?"

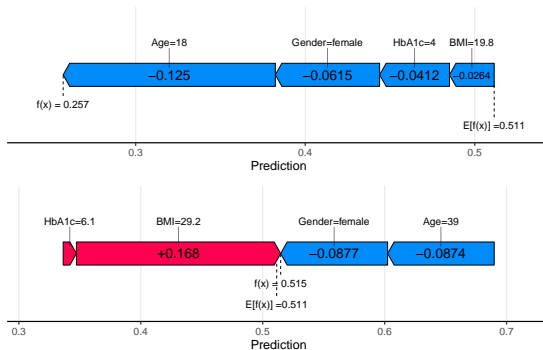
$$\text{GradSHAP}^7: \mathbb{E}_{\tilde{x} \sim X, \alpha \sim \mathcal{U}[0,1]} \left[(x_i - \tilde{x}_i) \frac{\partial f(\tilde{x} + \alpha(x - \tilde{x}))}{\partial x_i} \right]$$

DeepSHAP⁸ (rescale or reveal-cancel)

8

Interpretation:

- Decomposition of $f(x) - \mathbb{E}[f(X)]$ in feature-wise effects
- What is the marginal contribution of this feature?
 - answers a more natural question
 - axiomatic and theoretical foundation





"What are the features' contributions for predicting $f(x)$ compared to $\mathbb{E}_X[f(X)]$?"

GradSHAP⁷: $\mathbb{E}_{\substack{\tilde{x} \sim X \\ \alpha \sim \mathcal{U}[0,1]}} \left[(x_i - \tilde{x}_i) \frac{\partial f(\tilde{x} + \alpha(x - \tilde{x}))}{\partial x_i} \right]$

DeepSHAP⁸ (rescale or reveal-cancel)

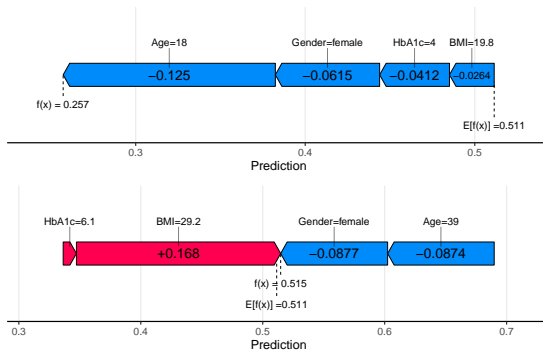
8

Interpretation:

- Decomposition of $f(x) - \mathbb{E}[f(X)]$ in feature-wise effects
- What is the marginal contribution of this feature?
 - answers a more natural question
 - axiomatic and theoretical foundation

Pitfalls:

- Still marginal (no conditional Shapley values)
- Higher computational costs
- Requires suitable (reference) dataset





"Is there a combined effect of features on $f(x)$?"

$$\text{Hessian: } \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$

ExpectedHessian⁹

$$\text{IntHessian}^9: (x_i - \tilde{x}_i)(x_j - \tilde{x}_j) \int_0^1 \int_0^1 \alpha \beta \frac{\partial^2 f(\tilde{x} + \alpha \beta (x - \tilde{x}))}{\partial x_i \partial x_j} d\alpha d\beta$$

9

$f(x) = 0.039$ vs. $f(x') = 0.569$

Interpretation:

- Decomposition of $f(x) - f(\tilde{x}) / f(x) - \mathbb{E}_X[f(X)]$ in feature-wise main (diagonal) and two-way interaction effects (**IntHessian**/ **ExpHessian**)
- Reveals local interaction effects and strength (**Hessian**)

Gender	0.353	-0.088	-0.069	-0.013
BMI		-0.505	0.152	0.036
Age			-0.334	0.051
HbA1c				-0.183
	Gender	BMI	Age	HbA1c



"Is there a combined effect of features on $f(x)$?"

$$\text{Hessian: } \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$

ExpectedHessian⁹

$$\text{IntHessian}^9: (x_i - \tilde{x}_i)(x_j - \tilde{x}_j) \int_0^1 \int_0^1 \alpha \beta \frac{\partial^2 f(\tilde{x} + \alpha\beta(x - \tilde{x}))}{\partial x_i \partial x_j} d\alpha d\beta$$

9

$f(x) = 0.039$ vs. $f(x') = 0.569$

Interpretation:

- Decomposition of $f(x) - f(\tilde{x}) / f(x) - \mathbb{E}_X[f(X)]$ in feature-wise main (diagonal) and two-way interaction effects (**IntHessian**/ **ExpHessian**)
- Reveals local interaction effects and strength (**Hessian**)

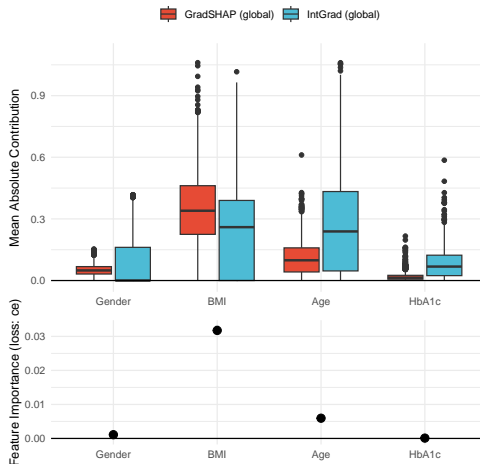
	Gender	BMI	Age	HbA1c
Gender	0.353	-0.088	-0.069	-0.013
BMI		-0.505	0.152	0.036
Age			-0.334	0.051
HbA1c				-0.183

Pitfalls:

- Higher computational costs
- Not possible for ReLU networks
 - vanishing 2nd derivative
- (similar to the other three groups)

- **So far:** Only explanations of individuals
- Aggregating local explanations for global insights (Lundberg et al. (2020))
- Gives relative importance among features
- Can we do so for other feature attribution methods?
 - Prediction-sensitive: Global feature's sensitivity (but depends on scaling)
 - Reference-based: Global effect against single reference
 - Shapley-based: Marginal global effect

⇒ Only explaining the model's prediction (not necessarily aligning with data-based or loss-based importance measures)



Key Takeaways:

- Feature Attribution Method \neq *feature attribution* \rightarrow answer different questions!
- Choice of method depends on the question: sensitivity, (baseline/marginal) attribution, or interaction
- Each reference value answers another question
- No recommendation for vague approximations (like LRP, Grad \times Input)
- Can be aggregated to global importance measures \rightarrow feature selection
- Prediction-based! \rightarrow "*What does the model see for the prediction*" – *not true to the data*

Future Work:

- Adopting methods for loss-based insights \rightarrow performance attribution
- Extending methods for conditional (not marginal) values

Thank you for your attention!



Slides, references and reproduction material



R Package `innsgit`

Contact

[Niklas Koenen](#)

Leibniz Institute for Prevention Research
and Epidemiology – BIPS

Achterstraße 30

D-28359 Bremen

koenen@leibniz-bips.de



I have no current or past relationships with commercial entities.