Shapley values and their uncertainties in complex regression models for epidemiological applications

Mark van de Wiel

Joint work with: Matteo Amestoy (Part I), Jeroen Hoogland (Part I); Giorgio Spadaccini (Part II), Marjolein Fokkema (Part II; U Leiden)

Epidemiology & Data Science, Amsterdam University Medical Centers, The Netherlands

Oslo, September 12, 2024

#### Part I: Regression models with many two-way interactions

- Linked shrinkage
- Shapleys plus intervals
- Example: canonical epidemiological problem

#### Part I: Regression models with many two-way interactions

- Linked shrinkage
- Shapleys plus intervals
- Example: canonical epidemiological problem

#### Part II: From regression to trees and back

- Translating trees to high-dimensional regression
- Shapleys plus intervals
- Example: academic performance of 1rst yr psychology students

## Part I: Setting

- Low-dimensional epidemiological cohort study (Helius)
  - Response: y, say cholesterol
  - Covariates: age, bmi, smoking, ethnicity, etc.
  - Our interest: n > p, say n = 1,000, p = 14
  - But large N available:  $N \approx 21.500$  (allows benchmarking)

## Part I: Setting

- Low-dimensional epidemiological cohort study (Helius)
  - Response: y, say cholesterol
  - Covariates: age, bmi, smoking, ethnicity, etc.
  - Our interest: n > p, say n = 1,000, p = 14
  - But large N available:  $N \approx 21.500$  (allows benchmarking)

#### Aims

- **1** Interpretable model that explains y
- Variable importance
- Ompetitive prediction

$$Y_i = \alpha + \sum_{j=1}^p \beta_j x_{ij} + \sum_{j,k:j \neq k} \beta_{jk} x_{ij} x_{ik} + \epsilon_i, \ \epsilon_i \sim N(0, \sigma^2)$$

$$Y_{i} = \alpha + \sum_{j=1}^{p} \beta_{j} x_{ij} + \sum_{j,k:j \neq k} \beta_{jk} x_{ij} x_{ik} + \epsilon_{i}, \ \epsilon_{i} \sim N(0, \sigma^{2})$$
$$\beta_{j} \sim N(0, \sigma^{2} \tau_{j}^{2})$$
$$\tau_{j} \sim C^{+}(0, 1)$$

 $\alpha \sim N(0, 10^2), \sigma^2 \sim IG(1, 0.001)$ 

$$Y_{i} = \alpha + \sum_{j=1}^{p} \beta_{j} x_{ij} + \sum_{j,k:j \neq k} \beta_{jk} x_{ij} x_{ik} + \epsilon_{i}, \ \epsilon_{i} \sim N(0, \sigma^{2})$$
$$\beta_{j} \sim N(0, \sigma^{2} \tau_{j}^{2}), \ \beta_{jk} \sim N(0, \sigma^{2} \tau_{j} \tau_{k} \tau_{int})$$
$$\tau_{j} \sim C^{+}(0, 1), \ \tau_{int} \sim U(0.01, 1)$$
$$\alpha \sim N(0, 10^{2}), \sigma^{2} \sim IG(1, 0.001)$$

$$Y_{i} = \alpha + \sum_{j=1}^{p} \beta_{j} x_{ij} + \sum_{j,k:j \neq k} \beta_{jk} x_{ij} x_{ik} + \epsilon_{i}, \ \epsilon_{i} \sim N(0, \sigma^{2})$$
$$\beta_{j} \sim N(0, \sigma^{2} \tau_{j}^{2}), \ \beta_{jk} \sim N(0, \sigma^{2} \tau_{j} \tau_{k} \tau_{int})$$
$$\tau_{j} \sim C^{+}(0, 1), \ \tau_{int} \sim U(0.01, 1)$$
$$\alpha \sim N(0, 10^{2}), \sigma^{2} \sim IG(1, 0.001)$$

Bayint: Coded in R-stan for flexibility

Models beats\* competitors like spike-and-slab, horseshoe, hierarchical lasso, adaptive lasso, etc on:

- Parameter estimation
- Prediction
- Variable selection
- ... and also RF (prediction)

<sup>\*</sup>See: https://arxiv.org/abs/2309.13998; vdW et al, *Epid Meth* 2024

# Comparison with RF (prediction)

#### Set-up

Outcomes: Cholesterol & SBP n = 1,000, p = 14 $n_{\text{test}} \approx 20,000$ 25 training-test splits **Competitors**:

- Bayint (our model)
- MainEff (our model with only main effects)
- OLS
- RF with default parameters
- RF with tuned parameters

# Comparison with RF (prediction)

#### Set-up

Outcomes: Cholesterol & SBP n = 1,000, p = 14  $n_{\text{test}} \approx 20,000$ 25 training-test splits **Competitors**:

- Bayint (our model)
- MainEff (our model with only main effects)
- OLS
- RF with default parameters
- RF with tuned parameters



Quantifying variable importance not trivial in regression model with interactions (Afshartous and Preston 2011)

Shapley values: metric quantifying variable importance per sample (Aas, Jullum, and Løland 2021)

Quantifying variable importance not trivial in regression model with interactions (Afshartous and Preston 2011)

Shapley values: metric quantifying variable importance per sample (Aas, Jullum, and Løland 2021)

- Originates from game theory
- Uniquely combines many desirable properties
- Applies to complex machine learners
- Expensive to compute, usually

#### Shapley, intuitive explanation (Molnar 2023)

- The variables enter a room in random order. All variables in the room (i.e. the coalition of players) participate in the game (= contribute to the prediction).
- Shapley value of a variable's realization: average change in the prediction when the variable joins the coalition
- Averaged over all possible coalitions (subsets).
- No refitting: predictions are marginalized w.r.t. non-players

Interventional Shapley: ignore dependencies between players and non-players

Explicit formula for the interventional Shapley for our model<sup>†</sup>: [e.g.  $x_{ij} = 49$ , age (j) for individual i]

<sup>&</sup>lt;sup>†</sup>See: https://arxiv.org/abs/2309.13998; vdW et al, *Epid Meth* 2024

Interventional Shapley: ignore dependencies between players and non-players

Explicit formula for the interventional Shapley for our model<sup>†</sup>: [e.g.  $x_{ij} = 49$ , age (j) for individual i]

$$\phi(x_{ij}) = \beta_j x_{ij} + \frac{1}{2} \left( \sum_{k:k \neq j} \beta_{jk} x_{ij} x_{ik} - \sum_{k:k \neq j} \beta_{jk} E[X_{ij} X_{ik}] \right),$$

<sup>&</sup>lt;sup>†</sup>See: https://arxiv.org/abs/2309.13998; vdW et al, *Epid Meth* 2024

### Shapleys + credible intervals

MCMC samples of  $\beta$  provides uncertainty estimate of  $\phi(x_{ij})$ 



### Intervals have good coverage



Figure: Coverages of 95% credible intervals for Shapley values of 200 random test individuals (quartiles shown). Estimated Shapley values are obtained from 500 random subsets of size n = 1,000.

### Distinguish contribution of main and interaction effects

Global importance:  $I_j = 1/n \sum_{i=1}^n |\phi(x_{ij})|$ .



## Comparison with hierarchical lasso

#### Hierarchical lasso (hlasso; Bien et al. 2013)

- Stronger link between interactions and main effects (selection)
- Computationally efficient
- Does not provide uncertainty estimates
- Shapleys hlasso: simply substitute parameter estimates into formula

#### Comparison

For variables Age, Etnicity:

- Stimate Shapley values for 1,000 random test individuals
- Obtain 'true' Shapley values from OLS estimates on Master set (N = 21500; p = 14; q = 85)
- Plot estimated vs true for 25 nearly non-overlapping training sets

### Shapley estimates (1: age)



Figure: Shapley values for 'age' over 1,000 random test individuals (dots) for 25 training sets (displays). X-axis: true Shapley values; Y-axis: estimated ones by Bayint (black) and hlasso (red).

### Shapley estimates (2: etn1)



Figure: Shapley values for 'etn1' over 1,000 random test individuals (dots) for 25 training sets (displays). X-axis: true Shapley values; Y-axis: estimated ones by Bayint (black) and hlasso (red).

Linked shrinkage useful concept for model with many interactions

Linked shrinkage useful concept for model with many interactions

Shapley values are useful for complex regression models

Linked shrinkage useful concept for model with many interactions

Shapley values are useful for complex regression models

Appropriate regularization allows good uncertainty quantification

Sometimes the world is non-linear and full of complex interactions

Then, tree learners (Random Forest, XGboost, etc) are a better alternative

Sometimes the world is non-linear and full of complex interactions

Then, tree learners (Random Forest, XGboost, etc) are a better alternative

Current practice in epidemiology:

- Use tree learners and Shap(ley) to select top k features
- Apply a linear (!!!) regression model to those features to perform inference

Sometimes the world is non-linear and full of complex interactions

Then, tree learners (Random Forest, XGboost, etc) are a better alternative

Current practice in epidemiology:

- Use tree learners and Shap(ley) to select top k features
- Apply a linear (!!!) regression model to those features to perform inference

**Main aim**: replace this inconsistent approach by one that is true to the tree learner

#### Develop an inferential approach true to the tree learner

- Translate tree learners to regression models
- ② Efficient computation of Shapley values for such models
- Inference: credible intervals for Shapley values
- Occomposition of these into linear contributions and remainder ones

#### Translate tree learners to regression models

- Ø Efficient computation of Shapley values for such models
- Inference on Shapley values: credible intervals
- Oecomposition of these into linear contributions and remainder ones

### From trees to regression: RuleFit

**RuleFit** (Friedman and Popescu 2008; Fokkema 2020): **linear regression** with **rule-based** interactions and non-linear terms.

Rules taken from Random Forest / Boosted Tree Ensemble  $x_2 < 3$   $r_1$   $x_2 \ge 3$   $r_2$   $r_3$   $r_4$  $r_4$ 

Rules + Linear terms  $\implies$  High interpretability Rules from Tree Ensemble  $\implies$  High performance

## RuleFit: Example

$$Y = X\beta = \hat{b}_1 x_1 + \hat{b}_1 x_2 + \hat{b}_3 x_3 + \hat{b}_4 x_4 + \hat{b}_5 x_5 + \hat{a}_1 I(x_2 < 3) + \hat{a}_2 I(x_2 \ge 3) + \hat{a}_3 I(x_2 < 3, x_5 < 7) + \hat{a}_4 I(x_2 < 3, x_5 \ge 7) + \cdots$$

$$X = \begin{bmatrix} 2 & -0.5 & 5 & 3.1 & 7.1 & 1 & 0 & 0 & 1 & \cdots \\ 1.5 & 3.7 & -0.1 & 4 & 2.4 & 0 & 1 & 0 & 0 & \cdots \\ 0.1 & 1.8 & 1 & 0 & 6 & 1 & 0 & 1 & 0 & \cdots \\ -0.9 & 4.1 & 2.9 & 2.2 & 2.1 & 0 & 1 & 0 & 0 & \cdots \\ -0.4 & 3 & 4.2 & 1.6 & -0.6 & 0 & 1 & 0 & 0 & \cdots \end{bmatrix}$$

Very many rules  $\Rightarrow$  Sparse model

- Translate tree learners to regression models
- In Efficient computation of Shapley values for such models
- Inference on Shapley values: credible intervals
- Oecomposition of these into linear contributions and rules

- Translate tree learners to regression models
- In Efficient computation of Shapley values for such models
- Inference on Shapley values: credible intervals
- Oecomposition of these into linear contributions and rules
- ightarrow Explicit formula for Shapley values  $\phi(x_{ij})$  based on RuleFit
  - Computationally competitive to TreeShap (Lundberg et al. 2020)
  - Advantage: achieves Aims 3 and 4 as well

### Decomposition and inference

- Translate tree learners to regression models
- Ø Efficient computation of Shapley values for such models
- Inference on Shapley values: credible intervals
- O Decomposition of these into linear contributions and rules

### Decomposition and inference

- Translate tree learners to regression models
- Ø Efficient computation of Shapley values for such models
- Inference on Shapley values: credible intervals
- O Decomposition of these into linear contributions and rules

#### Key ingredients for Aims 3 and 4:

- Better inference: Lasso (RuleFit)  $\rightarrow$  Horseshoe (Nalenz et al, 2018)
- Decomposition:
  - Rule generation: fit trees on residuals from linear model
  - $\bullet\,$  Joint fitting: many more rules than linear terms  $\rightarrow\,$  differential regularization

• 
$$\phi(x_{ij}) = \phi^{\mathsf{lin}}(x_{ij}) + \phi^{\mathsf{rule}}(x_{ij})$$

- Fit linear model
- It random forest<sup>‡</sup> on residuals (tweaked to counter over-fitting)
- Collect rules
- Sit horseshoe regression with unpenalized linear effects
- Obtain posteriors from MCMC sample
- Ompute Shapley values and uncertainties with formula

#### <sup>‡</sup>Could be replaced by one's favorite tree algorithm

#### Academic achievements; 1rst year psychology students

- Response: nr of credit points or mean grade
- Aim: predict and explain
- n = 638
- Covariates: gender, age, nationality (4), online\_test, test\_language, score\_Math, score\_Engl, score\_Psych, program
- Added five independent noise variables as negative controls
- Data available from pre package (Fokkema 2020)

Evaluation: cross-validated  $R^2$  (the higher, the better)

	Outcome	
Method	Credit Points	Mean Grade
RF	0.327	0.239
RuleFit	0.313	0.218
Tree	0.276	0.189
LassoReg	0.313	0.195
Ours§	0.320	0.260

§No good name yet!

### Predicting credit points: Shapleys + credible intervals



 $\text{Decomposition } j = \text{Score}\_\text{Engl:} \ (\overline{|\phi_j^{\text{lin}}|}, \overline{|\phi_j^{\text{rule}}|}) = (0.062, 0.016)$ 

### Predicting mean grade: Shapleys + credible intervals



Decomposition  $j = \text{Score}_{-}\text{Engl:} (\overline{|\phi_j^{\text{lin}}|}, \overline{|\phi_j^{\text{rule}}|}) = (0.003, 0.092)$ 

Translation to regression enables interpretation true to the tree-learner

Translation to regression enables interpretation true to the tree-learner

Shapley values are useful for complex regression models derived from tree-learners

Translation to regression enables interpretation true to the tree-learner

Shapley values are useful for complex regression models derived from tree-learners

Appropriate regularization allows uncertainty quantification but more work needed to show appropriate coverage of the intervals

#### References

Aas, Kjersti, Martin Jullum, and Anders Løland (2021). "Explaining individual predictions when features are dependent: More accurate approximations to Shapley values". In: Artificial Intelligence 298, p. 103502.

Afshartous, D. and R. A Preston (2011). "Key results of interaction models with centering". In: *Journal of Statistics Education* 19.3.

- Bien, Jacob, Jonathan Taylor, and Robert Tibshirani (2013). "A lasso for hierarchical interactions". In: *Annals of statistics* 41.3, p. 1111.
- Fokkema, Marjolein (2020). "Fitting Prediction Rule Ensembles with R Package pre". In: *Journal of Statistical Software* 92, pp. 1–30.
- Friedman, Jerome H and Bogdan E Popescu (2008). "Predictive Learning via Rule Ensembles". In: *The Annals of Applied Statistics*, pp. 916–954.
- Lundberg, Scott M et al. (2020). "From local explanations to global understanding with explainable AI for trees". In: *Nature machine intelligence* 2.1, pp. 56–67.
- Nalenz, Malte and Mattias Villani (2018). "Tree ensembles with rule structured horseshoe regularization". In: *The Annals of Applied Statistics* 12.4, pp. 2379–2408.
- van de Wiel, Mark A, Matteo Amestoy, and Jeroen Hoogland (2024). "Linked shrinkage to improve estimation of interaction effects in regression models". In: *Epidemiologic Methods* 13.1, p. 20230039.