

# OPTIMAL RULE-FIT ALGORITHM (ORFA)

Machine Learning Under an Optimization Lens

*Ryan Lucas & Paul Roeseler*

MIT SLOAN SCHOOL  
OF MANAGEMENT



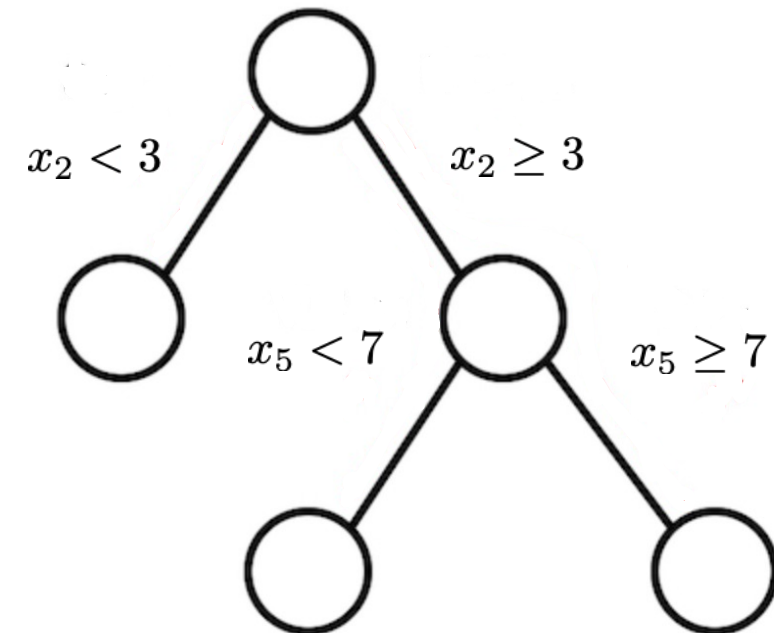
# MOTIVATION



## Decision trees and linear models uncover different types of effects

### Decision trees

- Uncover interaction effects



### Linear models

- Uncover linear relationships

$$\hat{Y} = X\hat{\beta}$$

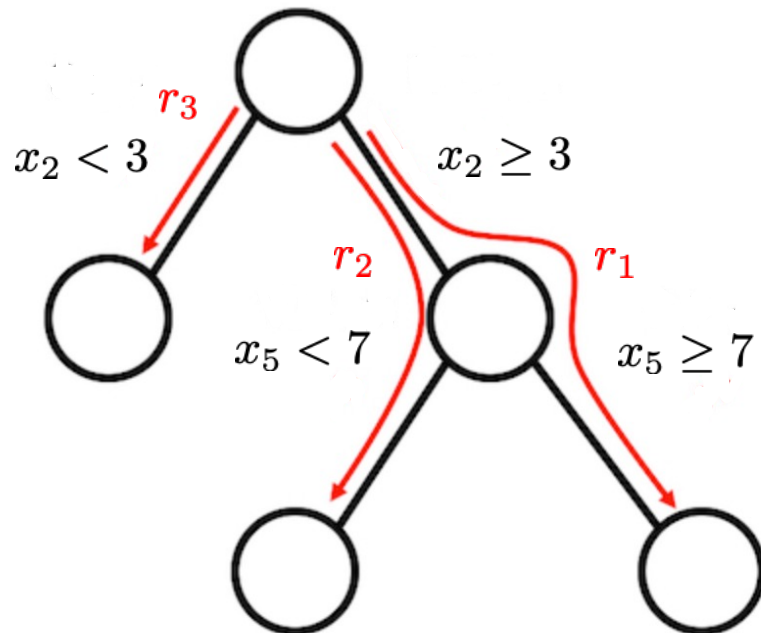
*But what if both types of effects are present?*

# MOTIVATION



## RuleFit algorithm (Friedman and Popescu, 2008)

Interpretable machine learning method using rules from a decision tree as features for a linear regression model



$$\hat{Y} = X\hat{\beta} + \hat{\delta}_1(\mathbb{1}\{x_2 \geq 3\} \cdot \{x_5 \geq 7\})$$

RuleFit adds rules as interaction features...

# MAJOR DRAWBACK



Greedy tree building methods (e.g., CART) require many splits to achieve strong performance – leads to great number of rules and overly sparse features

Few/Short Rules

Many/Long Rules



**Trade-off**



Interpretability

High Performance

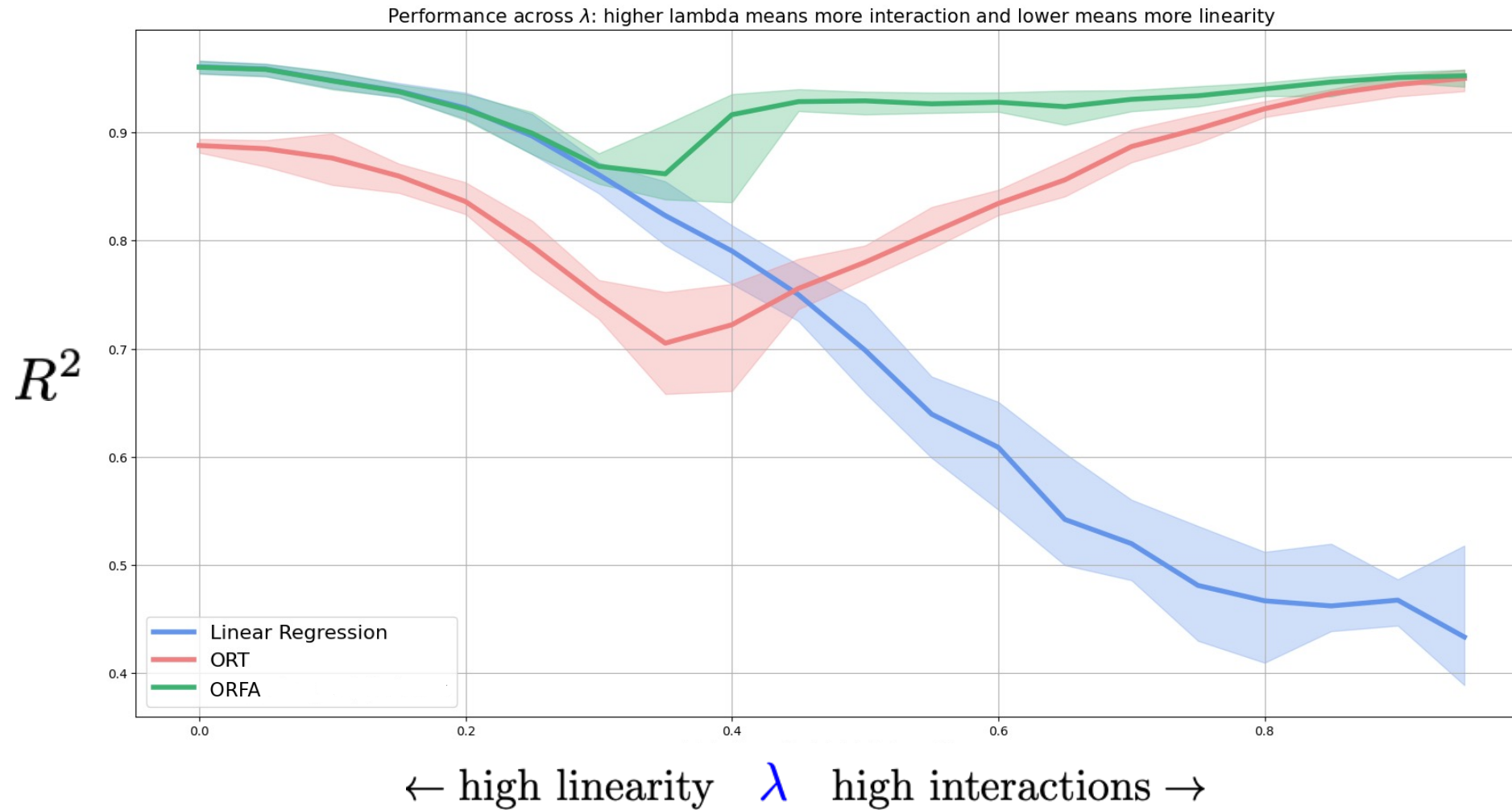
**Optimal Regression Trees (ORTs)** uncover true interaction effects in efficient number of splits and require only a single tree, resulting in fewer, more interpretable rules.

**We propose An Optimal RuleFit Algorithm (ORFA), combining ORTs and Linear Regression in a similar fashion to RuleFit**

# SIMULATIONS

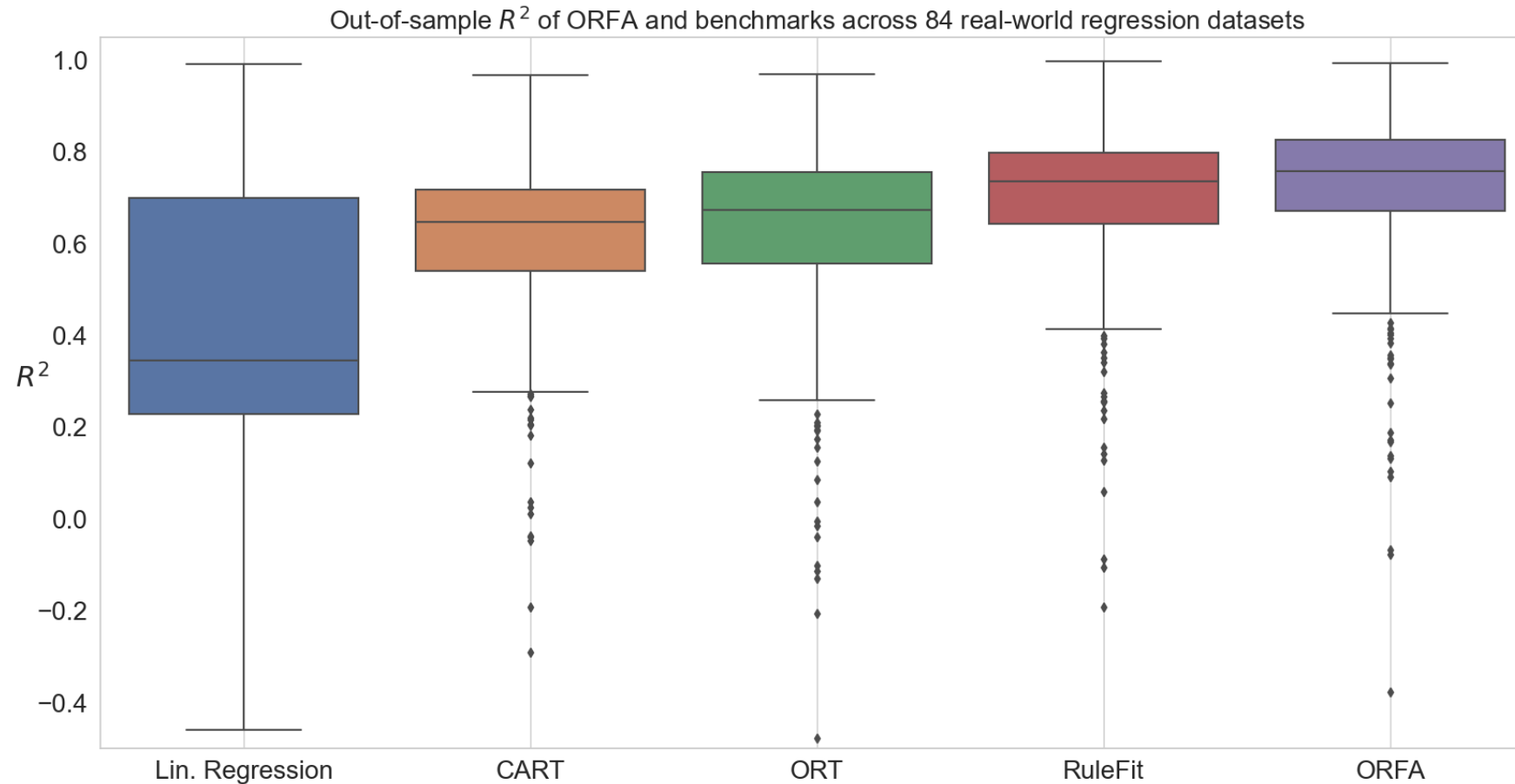


$$y_3 = \underbrace{\lambda \times \mathbb{1}\{x_3 \leq 0.3\} \times \mathbb{1}\{x_4 \geq -0.5\}}_{\text{interaction terms}} + \underbrace{(1 - \lambda) \times 0.5x_1 + 0.1x_2}_{\text{linear terms}} + \varepsilon$$



# BENCHMARK

Across 84 real-world regression datasets, provided by [PLMB](#), ORFA consistently ranks among the best methods



# BENCHMARK

Across 84 real-world regression datasets, provided by [PLMB](#), ORFA consistently ranks among the best methods



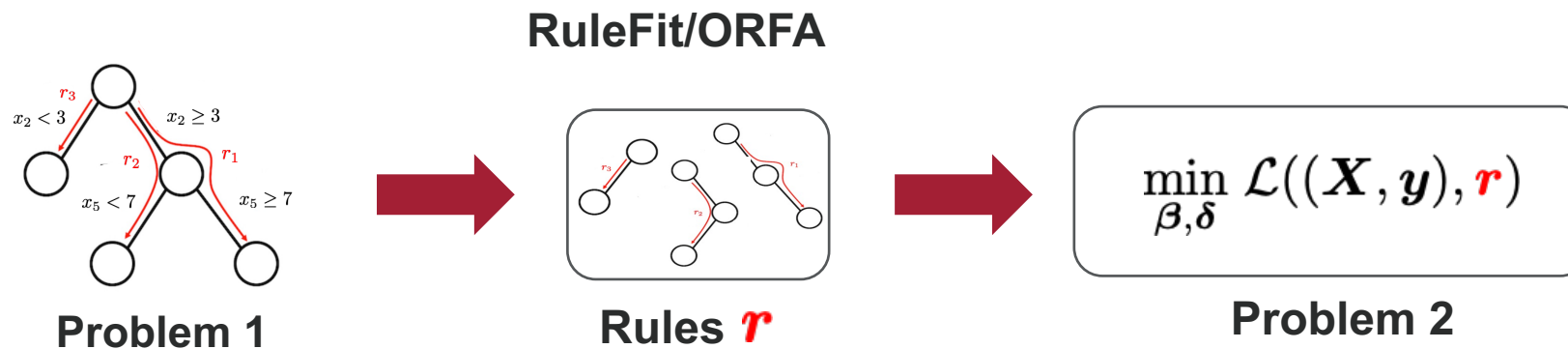
		Lin. Regression	CART	ORT	RuleFit	ORFA
Dataset	CPU	0.721	0.951	0.956	<b>0.976</b>	<b>0.978</b>
	Automobile	0.759	0.847	<b>0.879</b>	0.806	<b>0.874</b>
	Rabe	0.984	0.882	0.882	<b>0.986</b>	<b>0.987</b>
	Puma	0.375	0.567	<b>0.601</b>	0.571	<b>0.608</b>
	PW	0.710	0.780	0.762	<b>0.820</b>	<b>0.822</b>
	Wind	<b>0.754</b>	0.663	0.667	<b>0.754</b>	0.753
	Sleep Apnea	0.193	<b>0.845</b>	0.852	0.836	<b>0.844</b>
	Bodyfat	<b>0.974</b>	0.944	0.946	<b>0.974</b>	0.973
	CPU Small	0.707	0.936	0.947	<b>0.963</b>	<b>0.969</b>
	FRI	0.265	0.580	<b>0.684</b>	0.614	<b>0.749</b>
	Chatfield	0.851	0.704	0.679	0.781	<b>0.750</b>
	Geyser	<b>0.800</b>	<b>0.775</b>	0.755	0.779	0.762
⋮	⋮	⋮	⋮	⋮	⋮	
	Average	0.424	0.625	0.633	<b>0.707</b>	<b>0.724</b>

Table 1: Out-of-sample  $R^2$  across 84 real-world regression datasets provided by PMLB. The best performer on each dataset is highlighted in **blue**, while **purple** denotes the second best.

# BEYOND A HEURISTIC



The ORFA training is **disaggregated**. Rules are fed to regression and a new problem is solved.



## IORFA

$$\min_{\beta, \delta, r} \mathcal{L}(X, y, \text{tree})$$

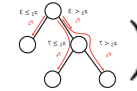
**Single Problem**



# INTEGRATED ORFA (IORFA)



Introducing IORFA, an integrated approach to solving  $\min_{\beta, \delta, \mathbf{r}} \mathcal{L}(\mathbf{X}, \mathbf{y}, \text{tree})$



IORFA is a modification to the MIO Formulation of ORT, introducing a regression objective:

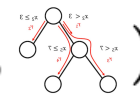
$$\min_{\beta, \delta} \sum_i (y_i - \mathbf{x}_i^T \beta - \sum_{t \in \mathcal{L}} \delta_t z_{i,t})^2 \quad z_{i,t} = \mathbb{1}\{x_i \in \text{leaf } t\}$$

subject to the usual constraints on  $\mathbf{z}$ ...

# INTEGRATED ORFA (IORFA)



Introducing IORFA, an integrated approach to solving  $\min_{\beta, \delta, \mathbf{r}} \mathcal{L}(\mathbf{X}, \mathbf{y}, \text{Diagram})$



IORFA is a modification to the MIO Formulation of ORT, introducing a regression objective:

$$\min_{\beta, \delta} \sum_i (y_i - \mathbf{x}_i^T \beta - \sum_{t \in \mathcal{L}} \delta_t z_{i,t})^2 \quad \text{with } z_{i,t} = \mathbb{1}\{x_i \in \text{leaf } t\}$$

subject to the usual constraints on  $\mathbf{z}$ ...

Think of  $\delta_t$  as fitting a coefficient on every group belonging to leaf nodes  $t \in \mathcal{L}$

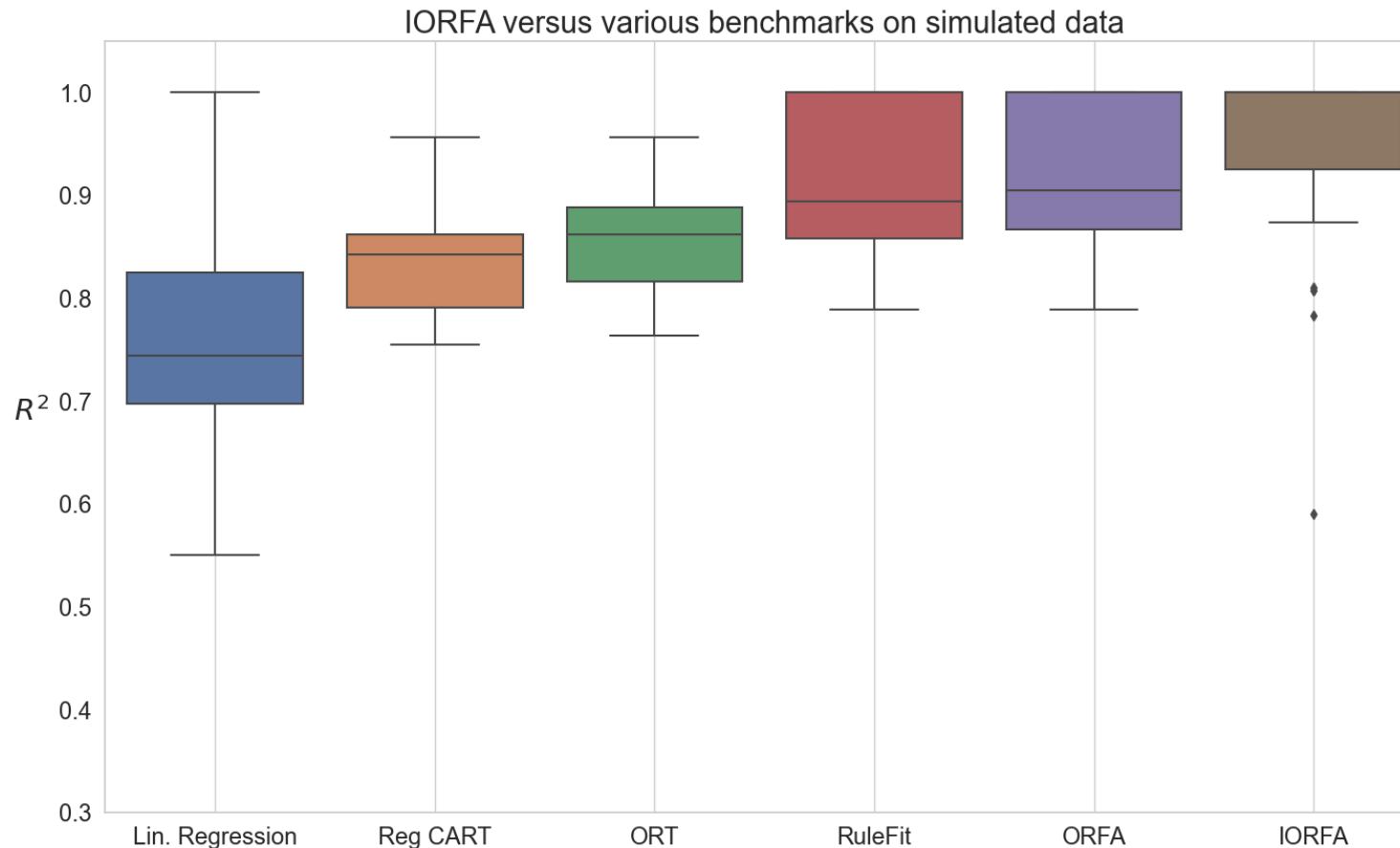
- Equivalent to fitting a parameter to every rule
- See [Appendix B](#) for the complete MIO formulation

# INTEGRATED ORFA (IORFA): RESULTS



Our resulting algorithm (IORFA) outperforms ORFA and RuleFit on a small simulated dataset

- We plan to extend these trials to real-world datasets in the coming weeks



# APPENDIX A: INTERPRETATION



Predicting bodyfat with one rule from OCT:

$$\text{Bodyfat}_i = 24.2 + 2.29 \cdot \text{Age}_i +, \dots, + 8.8 \cdot (\mathbb{1}\{\text{Weight}_i < 183.77\} \cdot \mathbb{1}\{\text{Age}_i < 37\} \cdot \mathbb{1}\{\text{Height}_i > 182.65\})$$

*If weight is less than 183.77 lbs, age is less than 38 and height is greater than 182.65 cm, then predicted bodyfat decreases by 8.8%, when all other feature values remain fixed.*

***This rule identifies a subgroup of tall, athletic young people with high weight but low bodyfat.***

# APPENDIX B: IORFA MIO FORMULATION



$$\begin{aligned}
 & \min \sum_{i=1}^n \theta_i \\
 & \text{s.t.} \quad \theta_i \geq (y_i - \mathbf{x}_i^T \boldsymbol{\beta} - \sum_{t \in \mathcal{L}} \delta_t z_{i,t}), \quad i = 1, \dots, n, \\
 & \quad \theta_i \geq -(y_i - \mathbf{x}_i^T \boldsymbol{\beta} - \sum_{t \in \mathcal{L}} \delta_t z_{i,t}), \quad i = 1, \dots, n, \\
 & \quad N_t = \sum_{i=1}^n z_{it}, \quad \forall t \in \mathcal{T}_L \\
 & \quad \mathbf{a}_m^\top \mathbf{x}_i \geq b_t - (1 - z_{it}), \quad i = 1, \dots, n, \quad \forall t \in \mathcal{T}_B, \quad \forall m \in A_R(t), \\
 & \quad \mathbf{a}_m^\top (\mathbf{x}_i + \boldsymbol{\epsilon}) \leq b_t + (1 + \epsilon_{\max})(1 - z_{it}), \quad i = 1, \dots, n, \forall t \in \mathcal{T}_B, \\
 & \quad \forall m \in A_L(t), \\
 & \quad \sum_{t \in \mathcal{T}_L} z_{it} = 1, \quad i = 1, \dots, n, \\
 & \quad z_{it} \leq l_t, \quad \forall t \in \mathcal{T}_L, \\
 & \quad \sum_{i=1}^n z_{it} \geq N_{\min} l_t, \quad \forall t \in \mathcal{T}_L, \\
 & \quad \sum_{j=1}^p a_{jt} = d_t, \quad \forall t \in \mathcal{T}_B, \\
 & \quad 0 \leq b_t \leq d_t, \quad \forall t \in \mathcal{T}_B, \\
 & \quad d_t \leq d_{p(t)}, \quad \forall t \in \mathcal{T}_B \setminus \{1\} \\
 & \quad z_{it}, l_t \in \{0, 1\}, \quad i = 1, \dots, n, \quad \forall t \in \mathcal{T}_L, \\
 & \quad a_{jt}, d_t \in \{0, 1\}, \quad j = 1, \dots, p, \quad \forall t \in \mathcal{T}_B
 \end{aligned}$$