

Ab Initio Nonparametric Variable Selection for Scalable Symbolic Regression with Large p

Shengbin Ye ^{1,2} Meng Li ¹¹Department Statistics, Rice University²Department of Statistics and Data Science, Northwestern University

Motivation

SR Algorithm		p	Has Irrelevant X
Transformer			
TPSR	(2023)	9	✗
RNN			
DySymNet	(2024)	9	✗
uDSR	(2022)	9	✗
DSR	(2021)	2	✗
Divide-and-conquer			
AlFeynman 2.0	(2020)	9	✗
Genetic Programming			
PySR	(2023)	6	✗
Operon	(2020)	5	✗
⋮		⋮	⋮

- Symbolic regression (SR) is **NP-hard**
- Most focus on **low-dimensional** problems (e.g., $p \leq 10$)
- Unrealistic settings: lack of irrelevant predictors
- No existing high-dimensional SR benchmark

PAN+SR: pre-screening framework for high-dimensional SR

Output y often depends on a **subset** $\mathcal{S}_0 \subseteq \{1, \dots, p\}$ of p_0 **relevant predictors**:

$$y = f(\mathbf{X}) = f(\mathbf{X}_{\mathcal{S}_0}),$$

where $p_0 = |\mathcal{S}_0| \ll p$.

Building on the **Parametrics Assisted by Nonparametrics (PAN)** framework in Ye et al. (JASA, 2024), we propose the PAN+SR framework.

Main Idea

- Nonparametric variable selection: $\mathbf{X} \mapsto \mathbf{X}_{\mathcal{S}}$
 - large $p \implies$ small p
 - SR search space $\downarrow\downarrow\downarrow$
- Perform SR on low-dimensional dataset $(y, \mathbf{X}_{\mathcal{S}})$

PAN Criterion

Step 1 **must select all** \mathcal{S}_0 : $\mathcal{S} \supseteq \mathcal{S}_0$

BART VIP Rank

We propose a novel Bayesian Additive Regression Tree (BART)-based variable selection method: **BART VIP Rank**

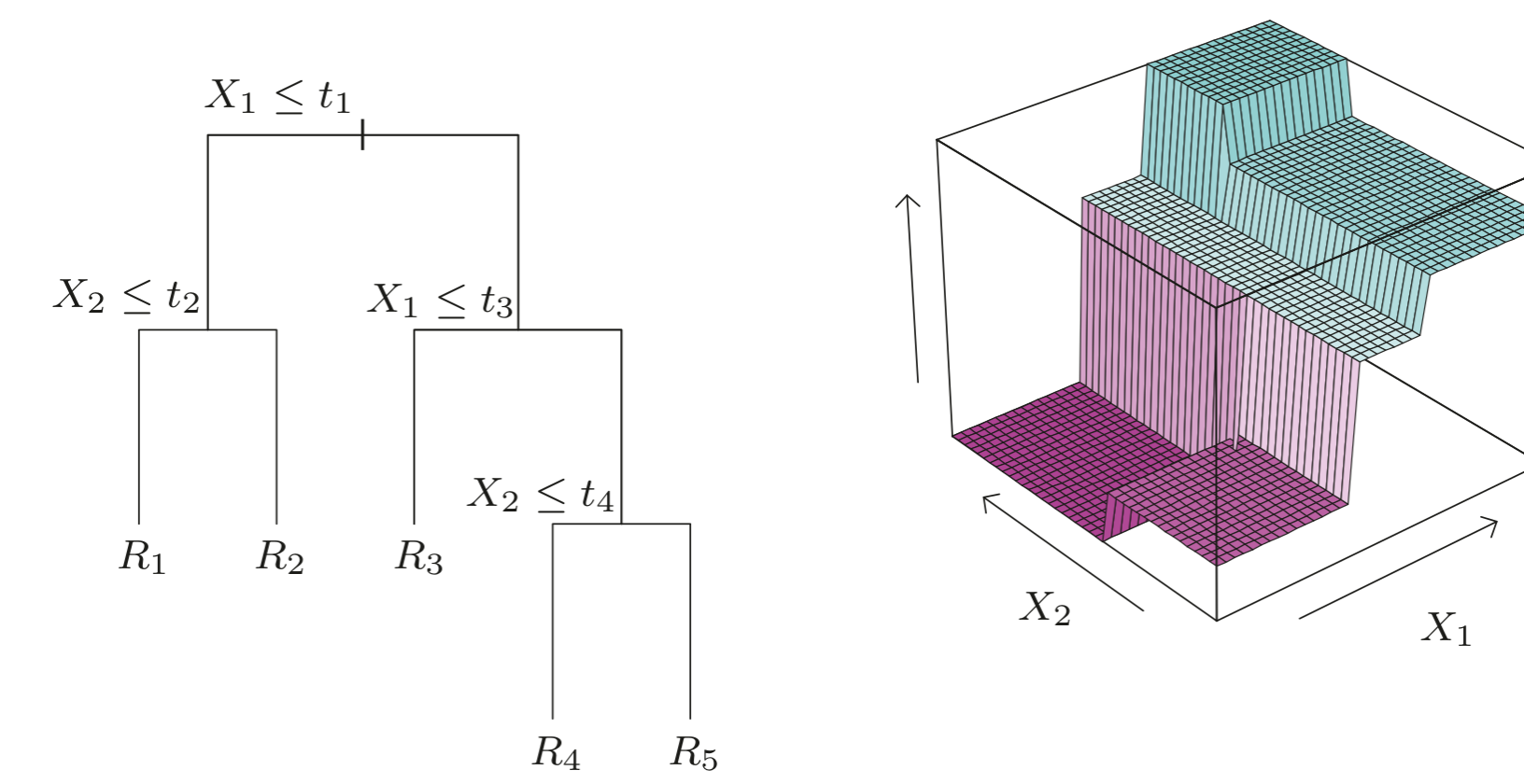


Figure 1. Visualization of BART.

Variable Inclusion Proportion (VIP)

A typical variable importance measure in BART is **variable inclusion proportion** (VIP):

$$q_j = \frac{1}{K} \sum_{k=1}^K \frac{c_{jk}}{c_{\cdot k}} \quad (\text{avg prop of splits on } x_j)$$

- Arbitrary scale**: how large is large?
- Tight range**: $0 \leq q_j \leq 1$
 - Small perturbation in threshold \implies different selections (sensitivity)

VIP Rank

Fit $L = 20$ independent BART models. Let $q_{j,\ell} = \text{VIP of } x_j \text{ in the } \ell\text{th fit}$, and let $R(q_{j,\ell}) = \text{ranking of } q_{j,\ell} \text{ within fit } \ell$. Define the **VIP Rank** for x_j as the average ranking over L model fits:

$$\bar{R}_j = \frac{1}{L} \sum_{\ell=1}^L R(q_{j,\ell}).$$

Under mild assumptions,

$$\bar{R}_j = \begin{cases} (1 + p_0)/2, & \text{if } x_j \text{ is relevant} \\ (p_0 + 1 + p)/2, & \text{otherwise} \end{cases}$$

Say $p_0 = 4$ and $p = 204$. Then, $\bar{R}_j = 2.5$ if x_j is relevant vs. $\bar{R}_j = 104.5$ otherwise.

Algorithm

- Fit $L = 20$ independent BART models on (y, \mathbf{X})
- Calculate BART VIP Rank $\bar{\mathbf{R}} = (\bar{R}_1, \dots, \bar{R}_p) \in \mathbb{R}^p$
- Apply Agglomerative Hierarchical Clustering on $\bar{\mathbf{R}}$
- Cut dendrogram to form 2 clusters: \mathcal{C}_{low} and $\mathcal{C}_{\text{high}}$
- Select x_j if $\bar{R}_j \in \mathcal{C}_{\text{low}}$

High-Dimensional SR Benchmark

We design a high-dimensional SR benchmark using 22 real-world datasets from PMLB and 100 synthetic datasets based on *Feynman Lectures on Physics*.

	SRBench	✦ New Benchmark
Irrelevant predictors \mathcal{S}_1	✗	✓
No. of predictors p	2 to 9	102 to 459
No. of relevant predictors p_0	2 to 9	2 to 9
Sample size n	100,000	500, 1000, 1500, 2000
Signal-to-noise ratio (SNR)	10 to noiseless	0.5 to noiseless
No. of trials per setting	10	10

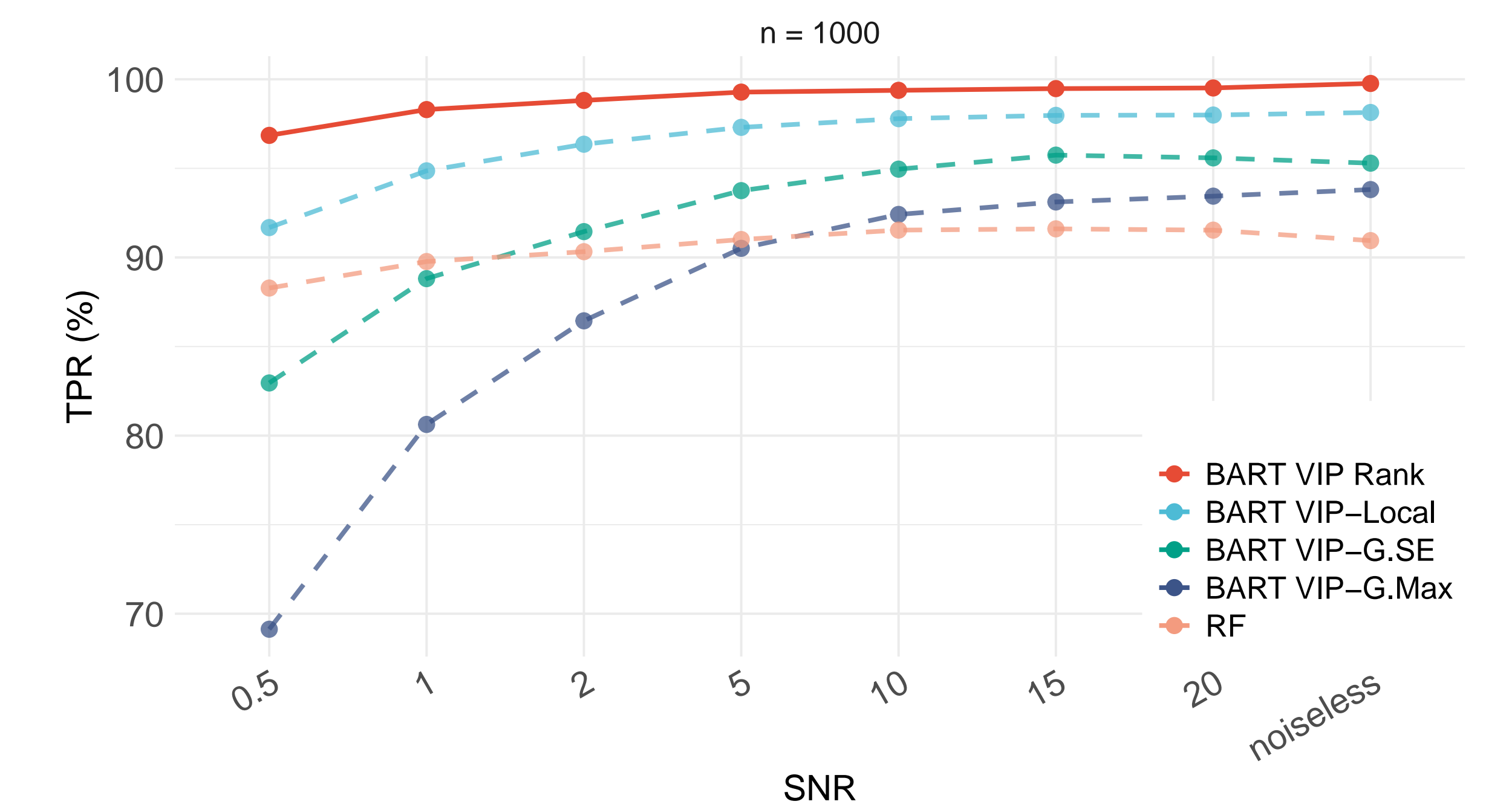


Figure 2. True positive rates on high-dimensional Feynman datasets.

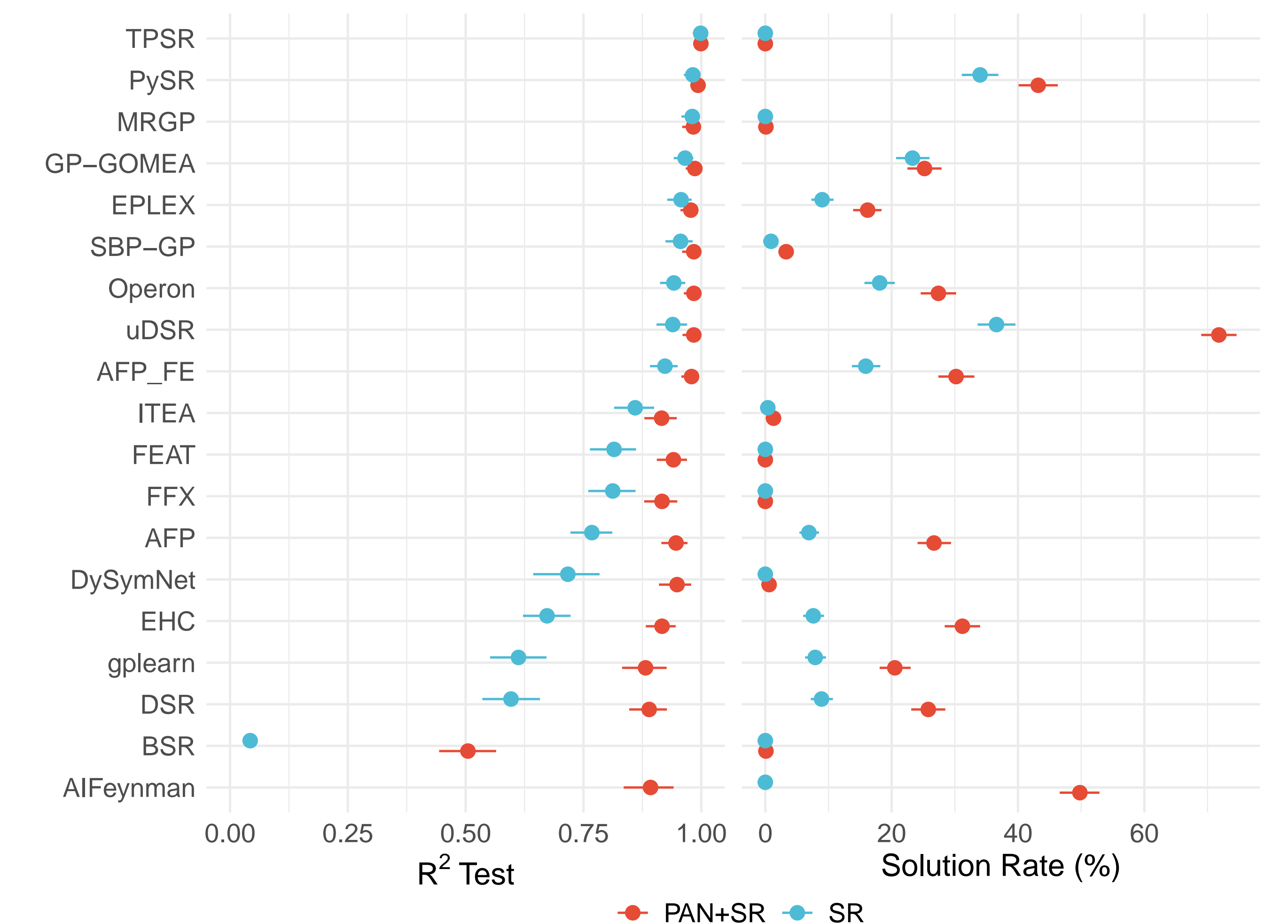


Figure 3. Performance of PAN+SR vs standalone SR on high-dimensional Feynman datasets.