

Can a new transparent algorithm predict better than its black-box counterparts?

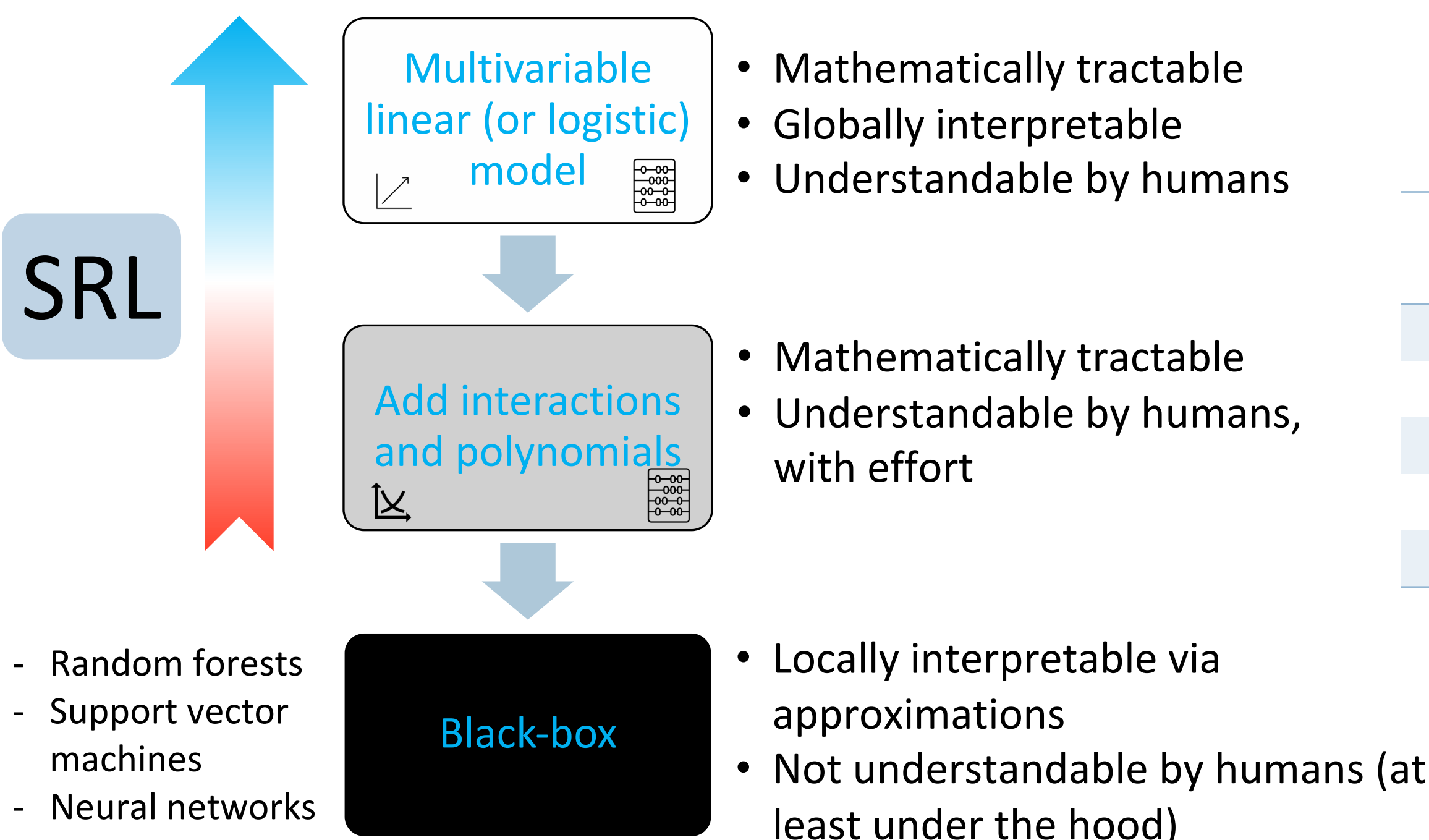
A benchmarking study for the Sparsity-Ranked Lasso using 112 diverse datasets

Ryan A. Peterson, PhD

Assistant Professor, Biostatistics & Informatics, University of Colorado Anschutz

TRANSPARENCY AND THE SPARSITY-RANKED LASSO

We developed¹ the **sparsity-ranked lasso (SRL)** as an alternative to black-box algorithms that *prefer transparency* in predictive models.

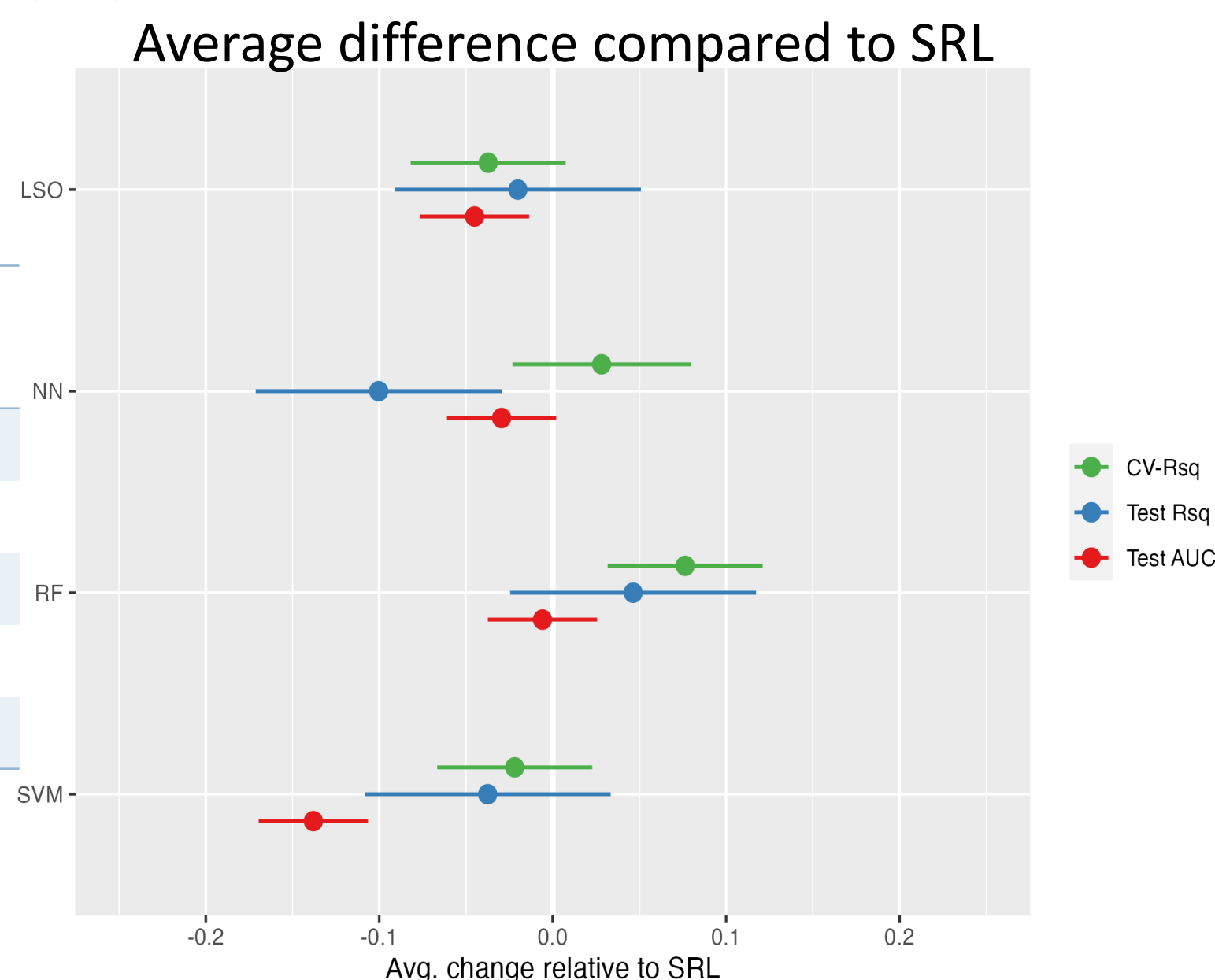


PREDICTIVE MODEL BAKEOFF

- Modeling methods, each with default settings used:
 - Transparent: lasso (LSO), sparsity-ranked lasso (SRL)
 - Black box: random forests (RF), support vector machines (SVMs), and neural networks (NN)²

Average performance across data sets:

	CV Rsq	Test Rsq	Test AUC
SRL	0.676	0.67	0.852
LSO	0.638	0.649	0.807
RF	0.752	0.716	0.846
SVM	0.654	0.632	0.714
NN	0.455	0.569	0.822



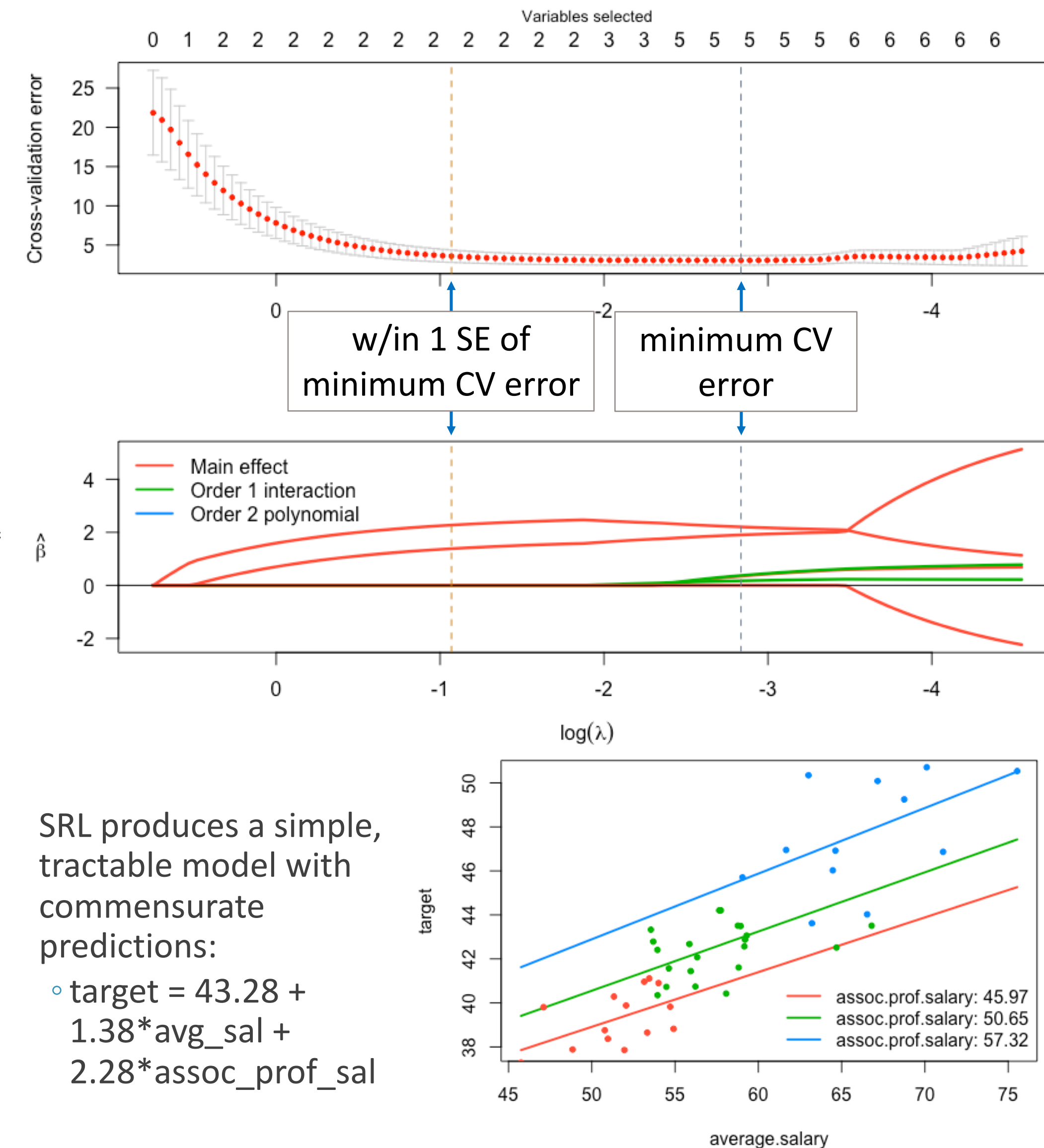
Compared to black-box methods, transparent ones had...

- best OOS R-squared in 32% of regression datasets
- best OOS AUC in 45% of classification datasets
- within 5% of best OOS R-squared/AUC in 70-80% datasets

CASE STUDY: "WIND 503" DATA SET (N=6574, P=14)

	SRL	Random forest
Tuning parameter values checked	101	3
Time to fit	4.12 seconds	~14min
Extra-sample R-squared	0.78	0.79
OOS R-squared	0.773	0.769

CASE STUDY: "FACULTY SALARIES" (N=50, P=4)

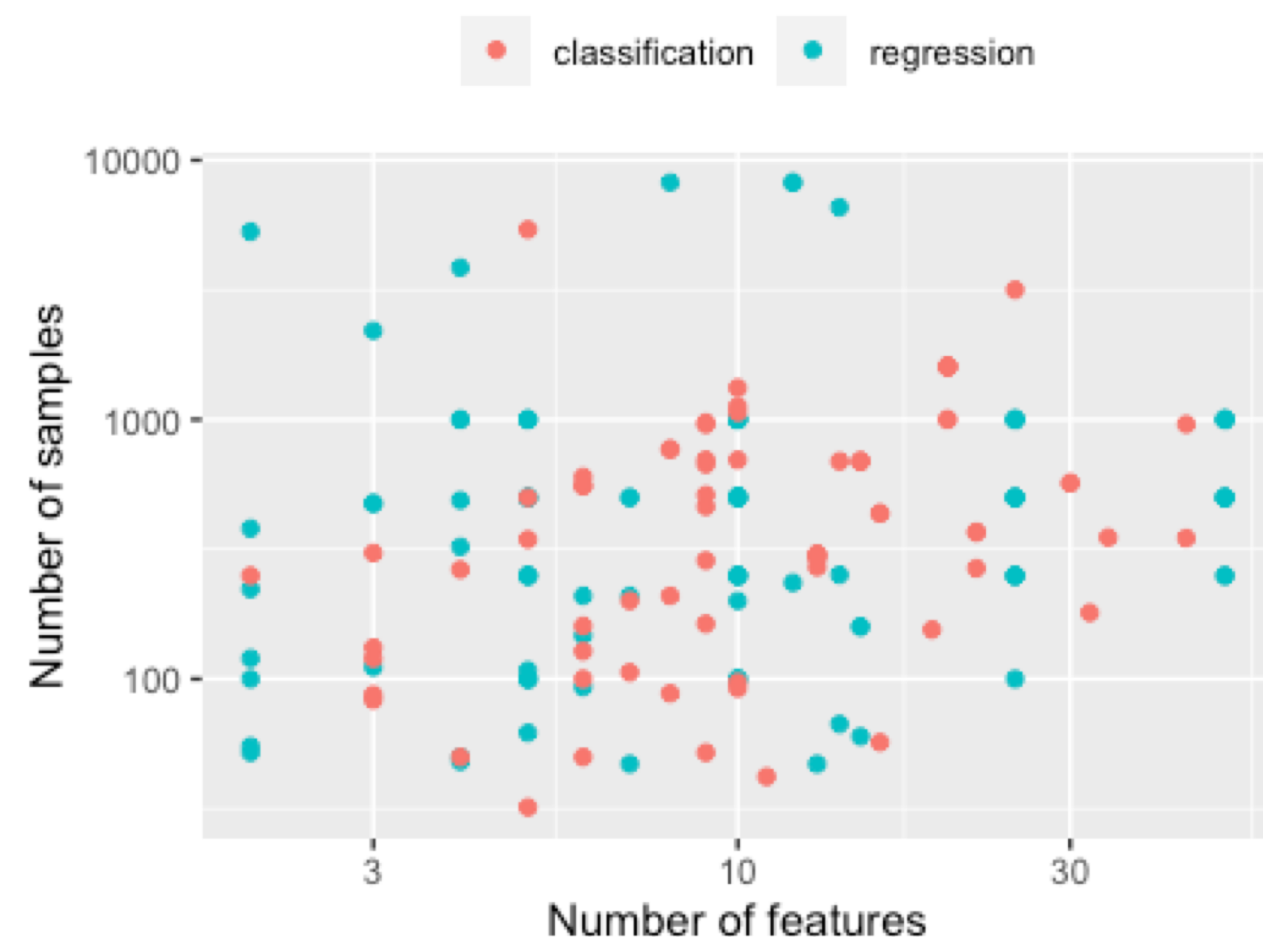


SRL produces a simple, tractable model with commensurate predictions:

$$\text{target} = 43.28 + 1.38 * \text{avg_sal} + 2.28 * \text{assoc_prof_sal}$$

A POPULATION OF DATA SETS

- N=112 datasets from the Penn Machine Learning Benchmarks Database
- A mix of simulated and real data sets, classification + regression problems
- Each data set split 75/25 into training/test set



CONCLUSIONS

- Our transparent algorithms sometimes predict better than black-box counterparts and most of the time perform comparably
 - At least for comparable data sets, e.g. not necessarily huge data sets.
- Takeaway: **always at least consider a transparent model.**

¹Peterson, R.A., Cavanaugh, J.E. Ranked sparsity: a cogent regularization framework for selecting and estimating feature interactions and polynomials. *ASA Advances in Statistical Analysis* (2022).
²Max Kuhn (2021). *caret: Classification and Regression Training*.