

SISSO: Selecting Sparsifying Operators from a Computational and Data Efficiency Perspective

AI³-2024, Paphos, Cyprus, 7 Nov 2024

Mario Boley^{1,2}

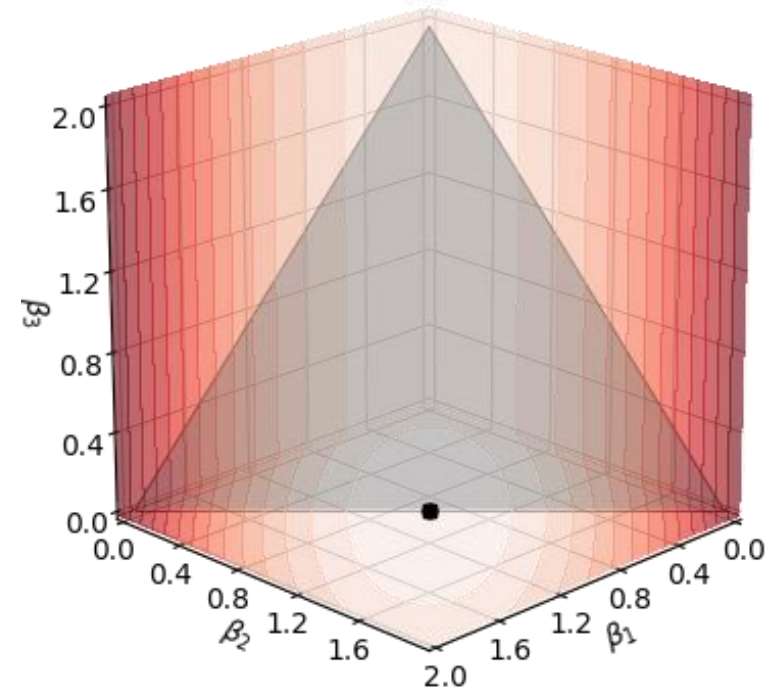
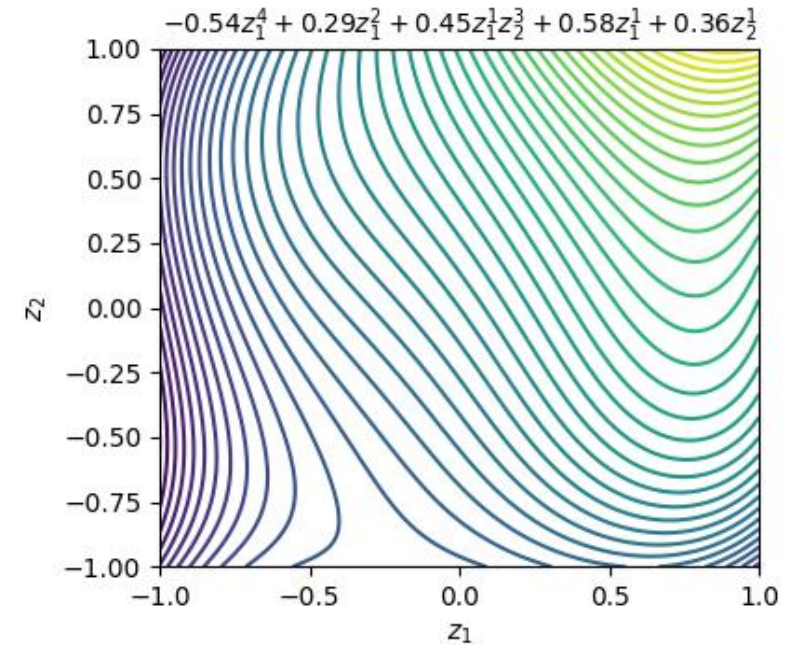
mboley@is.haifa.ac.il

Simon Teshuva¹, Felix Luong¹, Daniel Schmidt¹, Lucas Foppa³, and Matthias Scheffler³

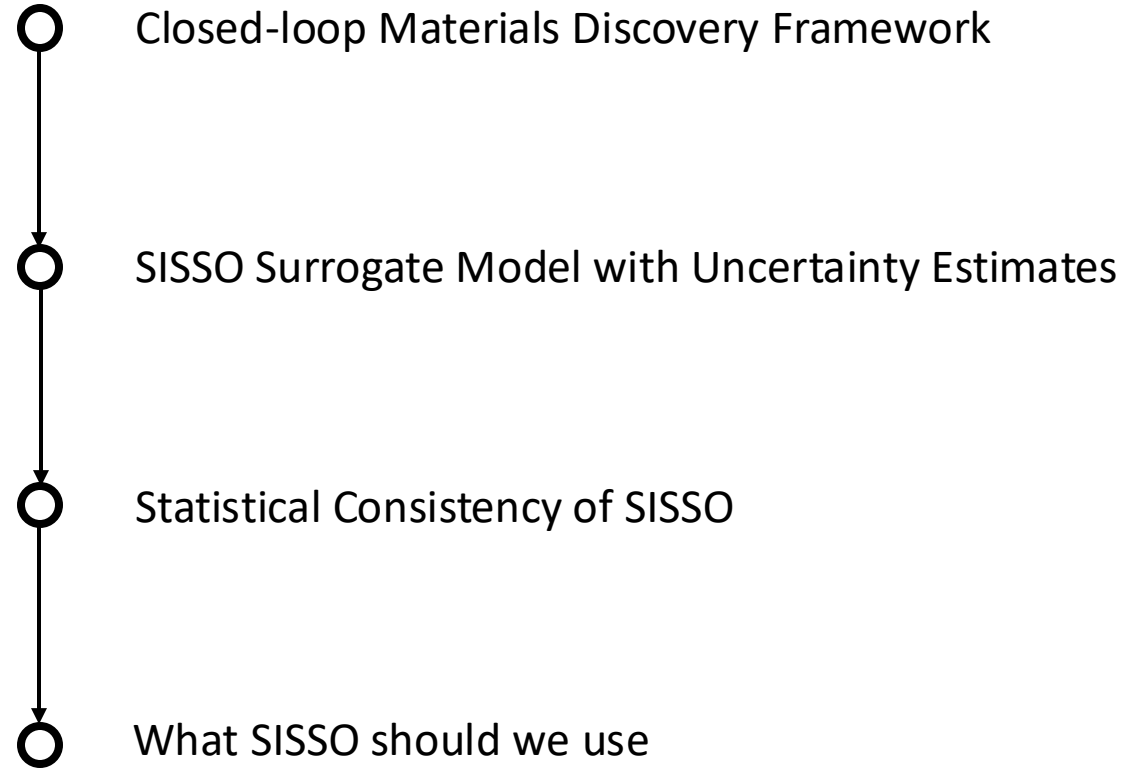
¹Department of Information Systems, University of Haifa

²Department of Data Science and AI, Monash University

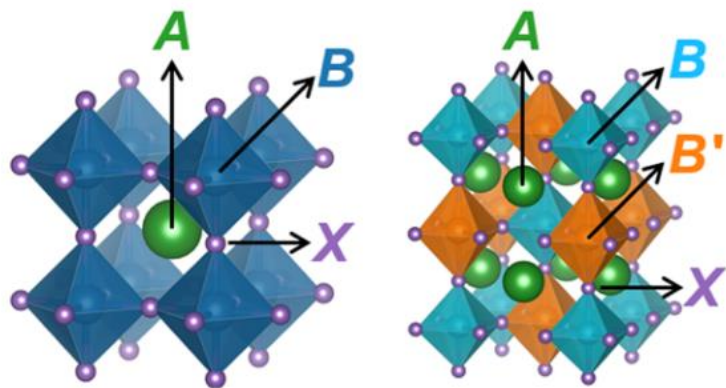
³NOMAD Laboratory at FHI of Max-Planck-Gesellschaft



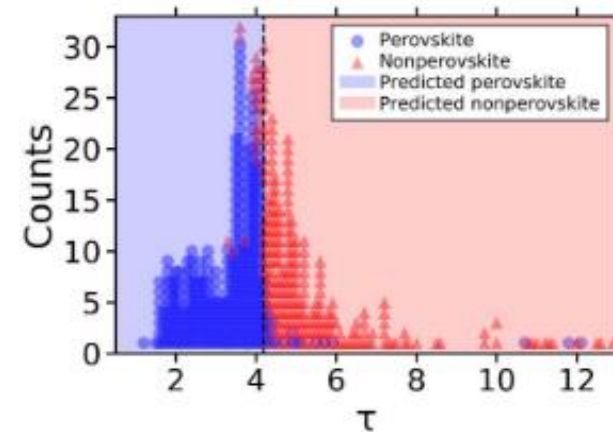
How it all connects



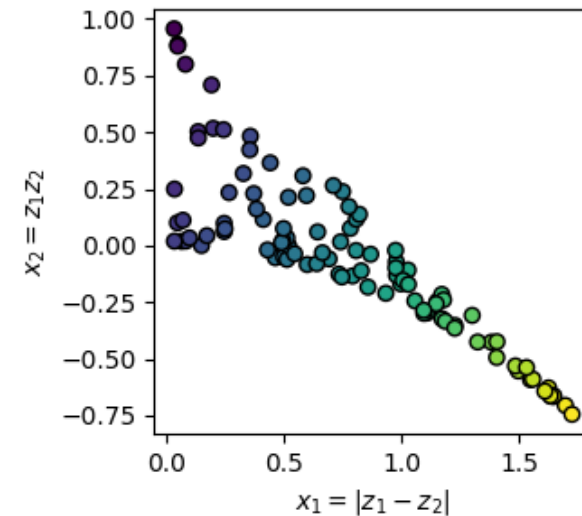
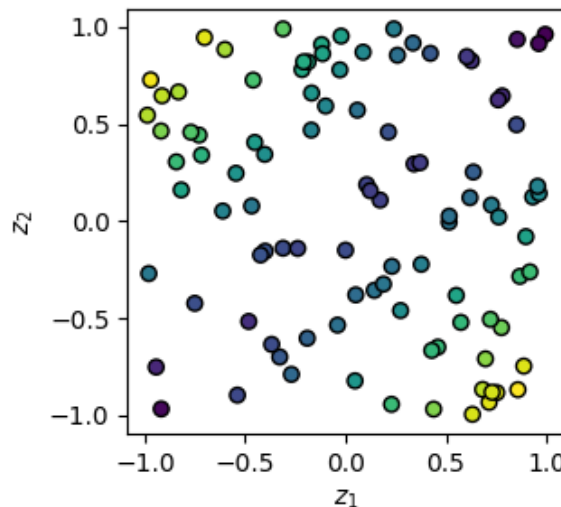
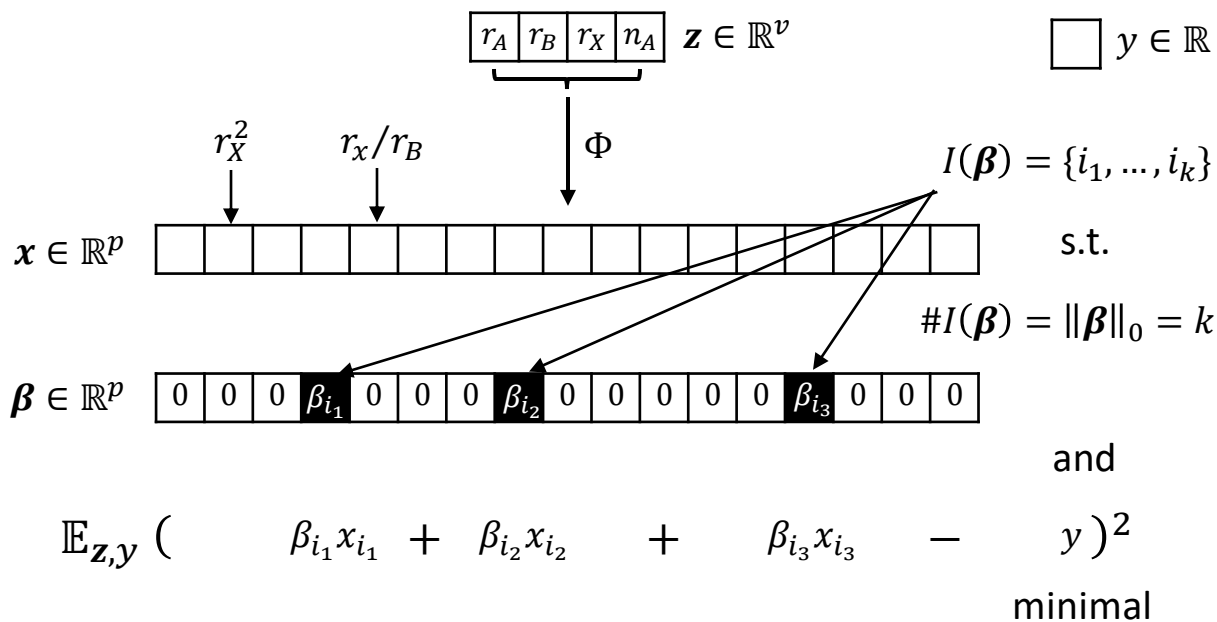
SISSO: Symbolic Regression for Materials Properties



$$\log \frac{P(\text{stable})}{1 - P(\text{stable})} = \beta_1 \frac{r_X}{r_B} + \beta_2 n_A^2 - \beta_3 \frac{n_A r_A / r_B}{\ln(r_A / r_B)}$$



[Bartel, C. J., et al. (2019). *New tolerance factor to predict the stability of perovskite oxides and halides*. *Sci. Adv.* 5(2).]



[Ouyang et al. (2018). *SISSO: A compressed-sensing method for low-dimensional descriptors*. *Phys. Rev. Mater.* 2(8)]

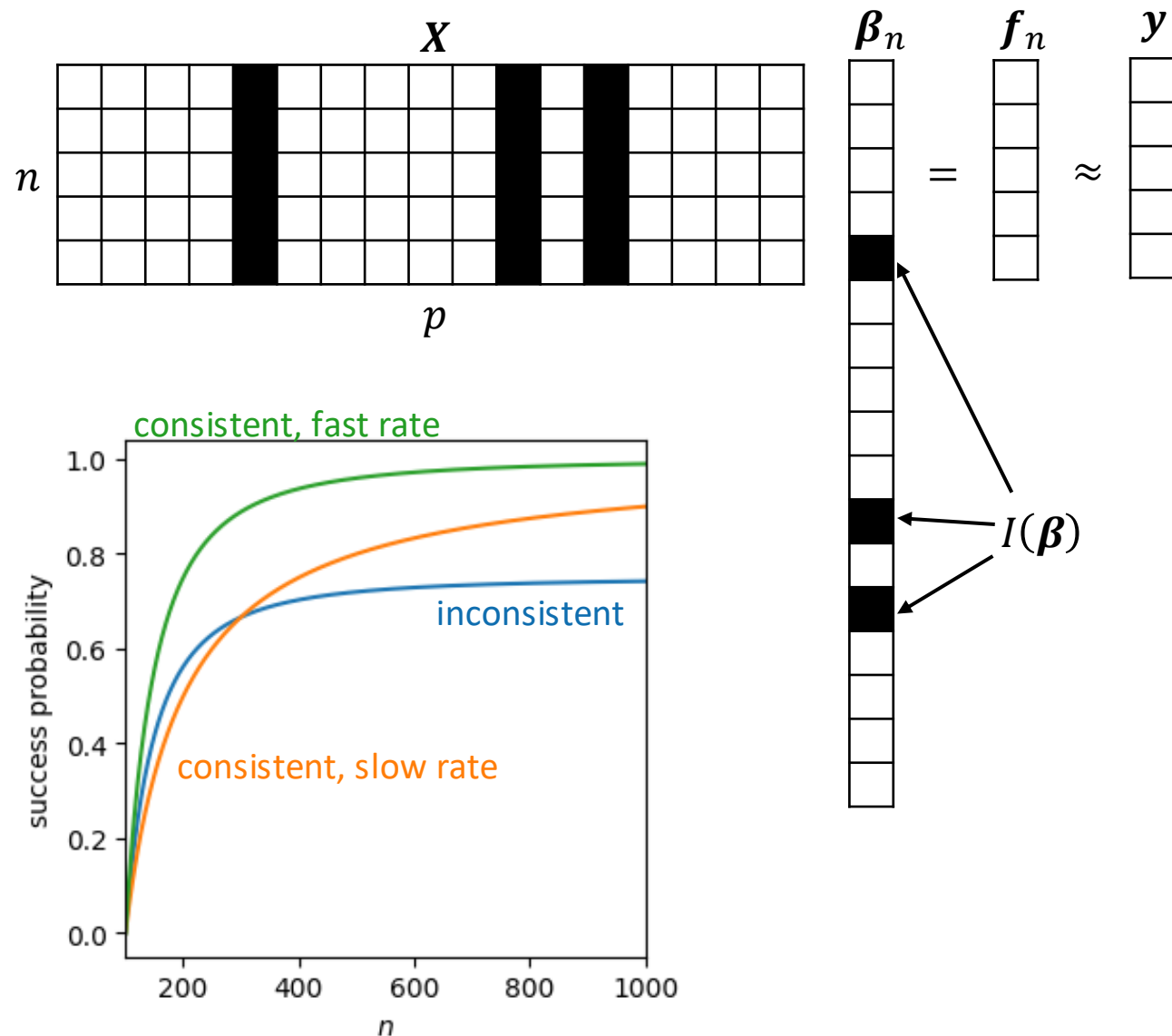
Need to Select Subset via Data Sample

Given:

- **input** matrix $X \in \mathbb{R}^{n \times p}$, **output** vector $y \in \mathbb{R}^n$ with rows sampled w.r.t. joint x, y distribution
- prescribed **sparsity**/complexity $k \in \mathbb{N}$
- typically assume $k < n \ll p$

Goal:

- identify $\beta_* = \operatorname{argmin}\{\mathbb{E}(y - x^T \beta)^2 : \#I(\beta) = k\}$
- via sparse estimate β_n with $\#I(\beta_n) = k$
- computationally **efficiently**, i.e., in time $O(knp)$
- **consistently**, i.e., $\lim_{n \rightarrow \infty} P(I(\beta_n) = I(\beta_*)) = 1$
- with as **fast a rate** as possible



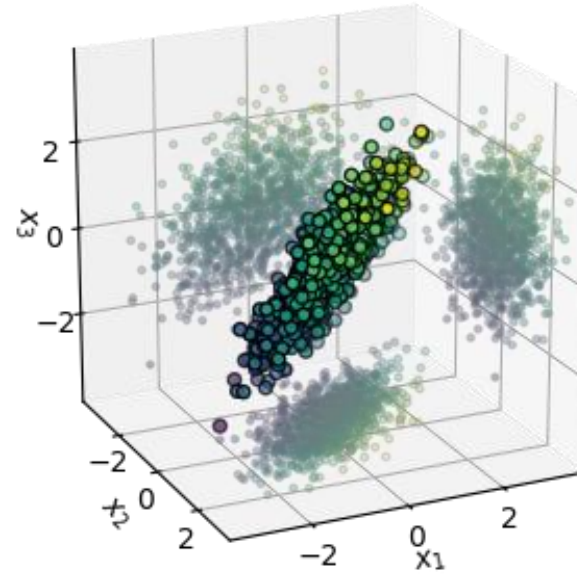
There are many methods... that fail to reach goals

Best-subset-search:

find $\beta_n^{\text{BSS}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}\beta\|^2: \#I(\beta) = k\}$

consistent (ordinary least squares parameter consistency)

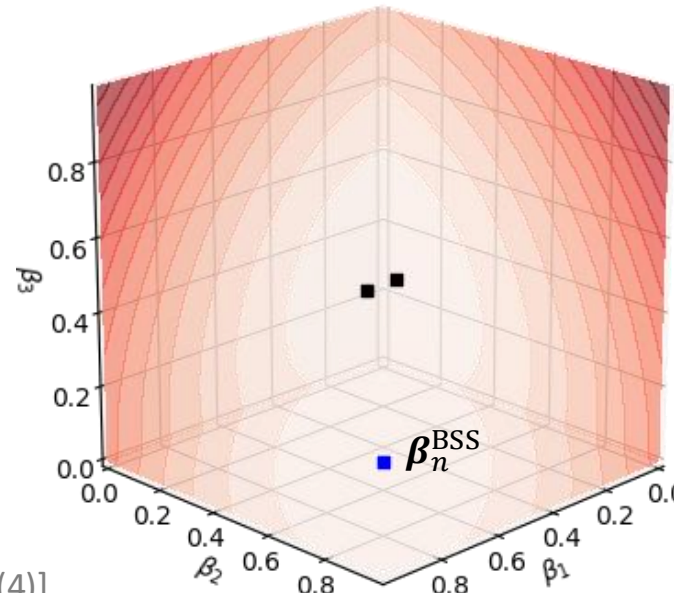
but computationally inefficient $O(C_{p,k}(nk^2 + k^3))$



$$y = 0.5x_1 + 0.5x_2$$

$$x \sim N_3(0, C)$$

$$C = \begin{bmatrix} 1 & -3/4 & 0.3 \\ -3/4 & 1 & 0.3 \\ 0.3 & 0.3 & 1 \end{bmatrix}$$



There are many methods... that fail to reach goals

Best-subset-search:

find $\beta_n^{\text{BSS}} = \operatorname{argmin}\{\|y - X\beta\|^2: \#I(\beta) = k\}$

consistent (ordinary least squares parameter consistency)

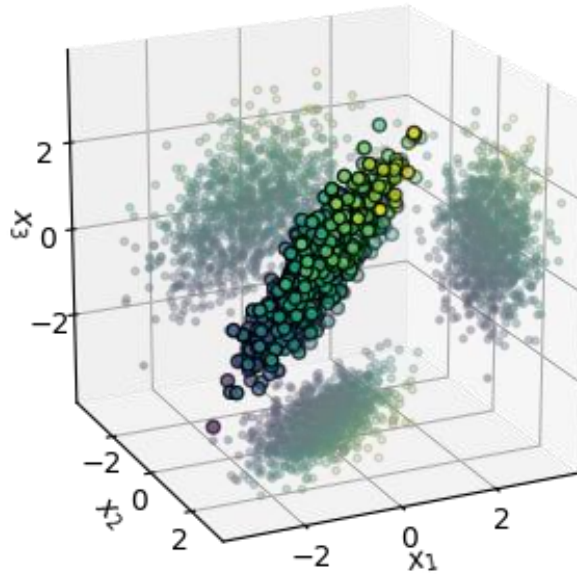
but computationally inefficient $O(C_{p,k}(nk^2 + k^3))$

LASSO:

find $\beta_n^{\text{LAS}} = \operatorname{argmin}\{\|y - X\beta\|^2: \|\beta\|_1 \leq c_k\}$

computationally efficient $O(knp + k^3)$

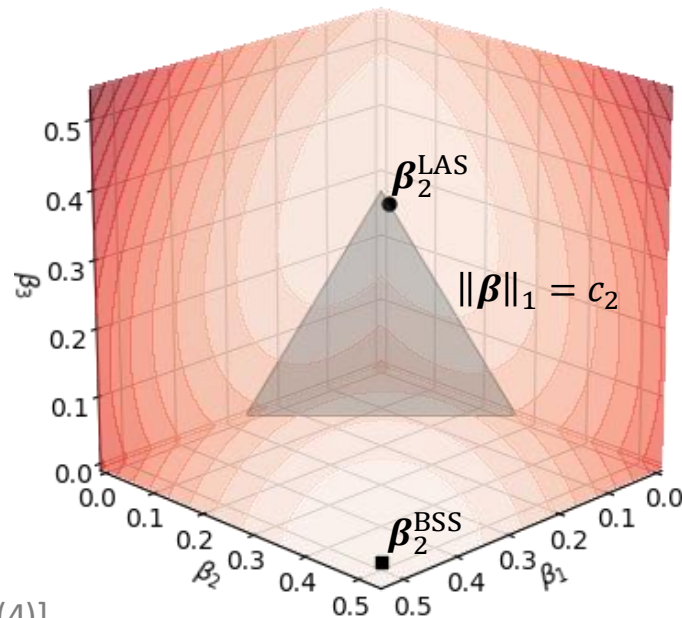
but inconsistent for non-trivial correlation structure



$$y = 0.5x_1 + 0.5x_2$$

$$x \sim N_3(0, C)$$

$$C = \begin{bmatrix} 1 & -3/4 & 0.3 \\ -3/4 & 1 & 0.3 \\ 0.3 & 0.3 & 1 \end{bmatrix}$$



There are many methods... that fail to reach goals

Best-subset-search:

$$\text{find } \boldsymbol{\beta}_n^{\text{BSS}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2: \#I(\boldsymbol{\beta}) = k\}$$

consistent (ordinary least squares parameter consistency)

but computationally **inefficient** $O(C_{p,k}(nk^2 + k^3))$

LASSO:

$$\text{find } \boldsymbol{\beta}_n^{\text{LAS}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2: \|\boldsymbol{\beta}\|_1 \leq c_k\}$$

computationally **efficient** $O(knp + k^3)$

but **inconsistent** for non-trivial correlation structure

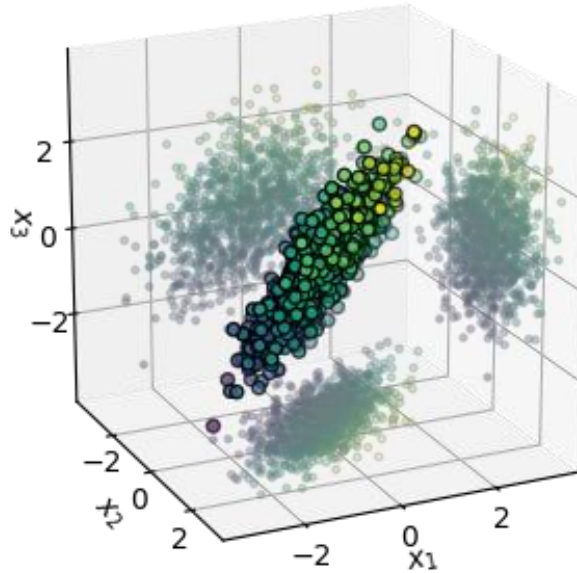
Thresholded Minimum-Norm Least Squares:

$$\text{find } \boldsymbol{\beta} = \operatorname{argmin}\left\{\lim_{\lambda \rightarrow 0_+} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_2^2\right\}$$

$$\text{and set } \beta_j^{\text{TLS}} = \begin{cases} \beta_j, & \text{if } |\beta_j| \text{ among } k \text{ largest} \\ 0, & \text{otherwise.} \end{cases}$$

consistent (although rate can be slow)

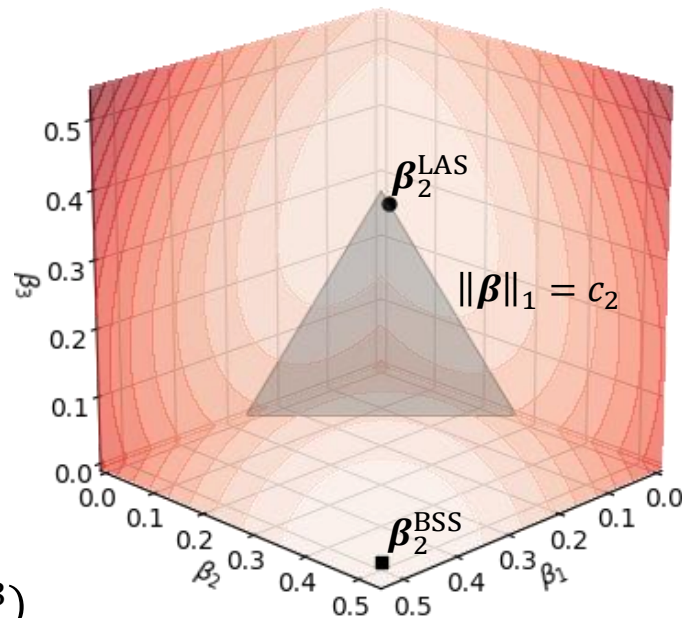
computationally **inefficient** $O(np^2 + p^3)$ or $O(n^2p + n^3)$



$$y = 0.5x_1 + 0.5x_2$$

$$x \sim N_3(0, C)$$

$$C = \begin{bmatrix} 1 & -3/4 & 0.3 \\ -3/4 & 1 & 0.3 \\ 0.3 & 0.3 & 1 \end{bmatrix}$$



There are many methods... that fail to reach goals

Best-subset-search:

find $\beta_n^{\text{BSS}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}\beta\|^2: \#I(\beta) = k\}$

consistent (ordinary least squares parameter consistency)

but computationally **inefficient** $O(C_{p,k}(nk^2 + k^3))$

LASSO:

find $\beta_n^{\text{LAS}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}\beta\|^2: \|\beta\|_1 \leq c_k\}$

computationally **efficient** $O(knp + k^3)$

but **inconsistent** for non-trivial correlation structure

Adaptive LASSO

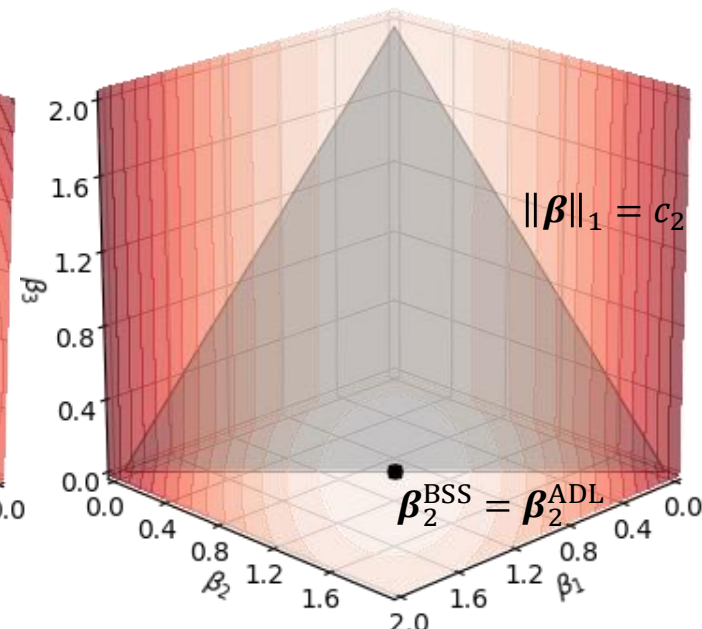
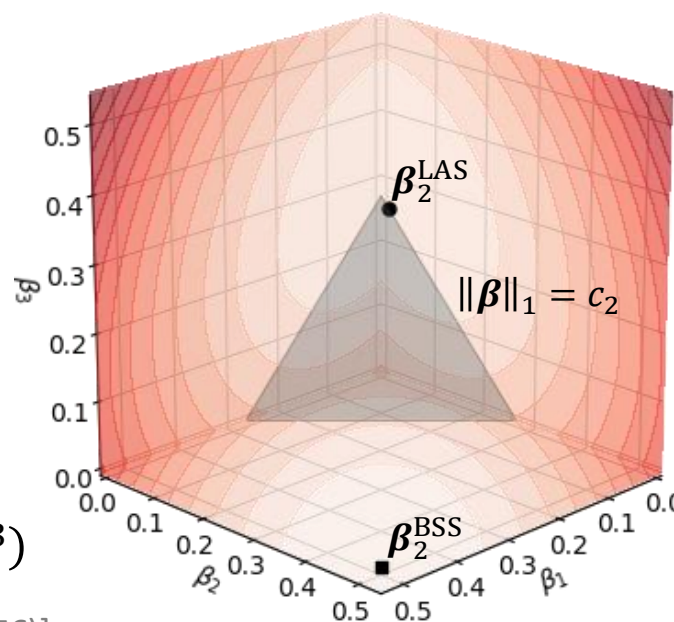
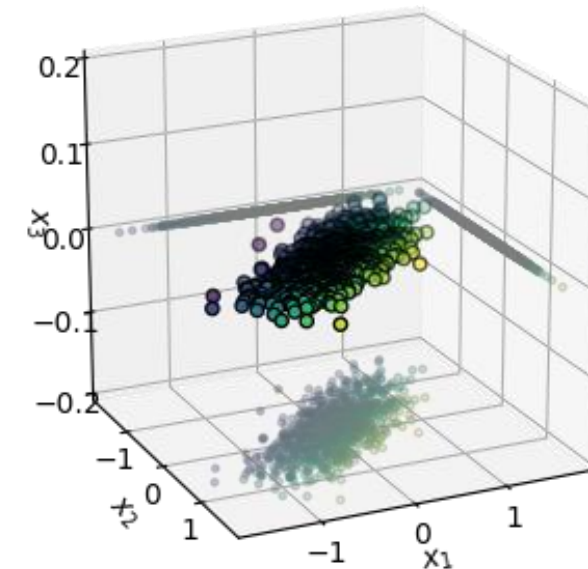
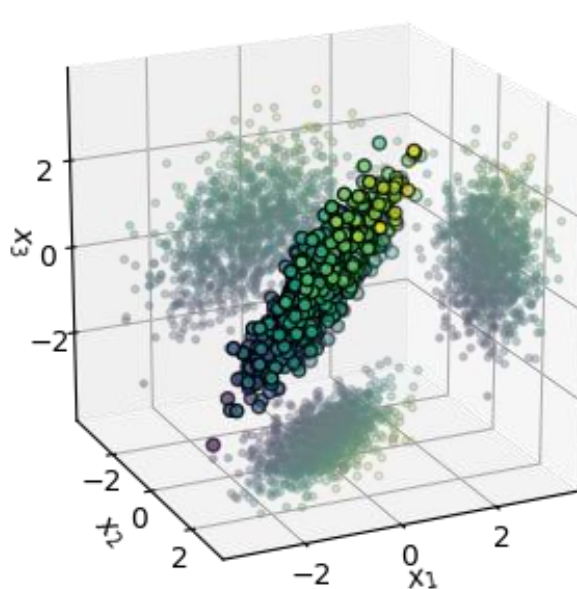
find $\alpha = \operatorname{argmin}\left\{\lim_{\lambda \rightarrow 0_+} \|\mathbf{y} - \mathbf{X}\alpha\|^2 + \lambda \|\alpha\|_2^2\right\}$

and $\beta' = \operatorname{argmin}\|\mathbf{y} - \mathbf{Z}\beta'\|^2 + \lambda_k \|\beta'\|_1$

and $\beta_j = |\alpha_j| \beta'_j$ where $z_{i,j} = |\alpha_j| x_{i,j}$

consistent (oracle rate in parameter reconstruction)

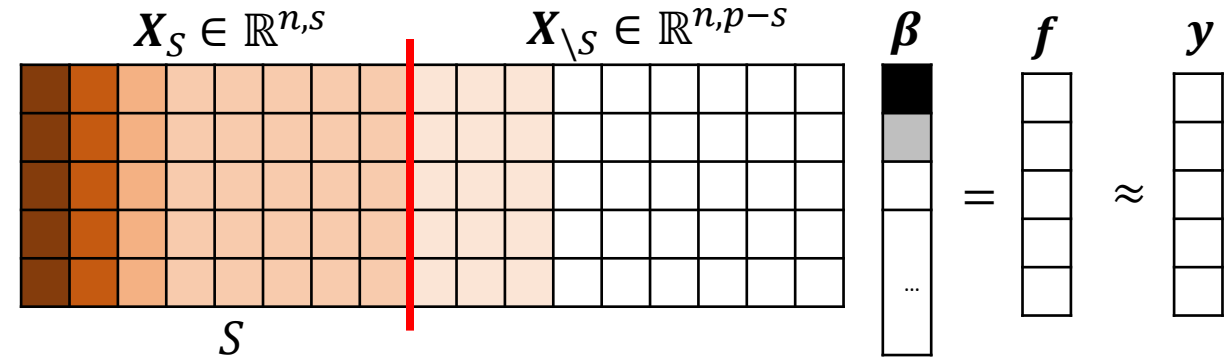
computationally **inefficient** $O(np^2 + p^3)$ or $O(n^2p + n^3)$



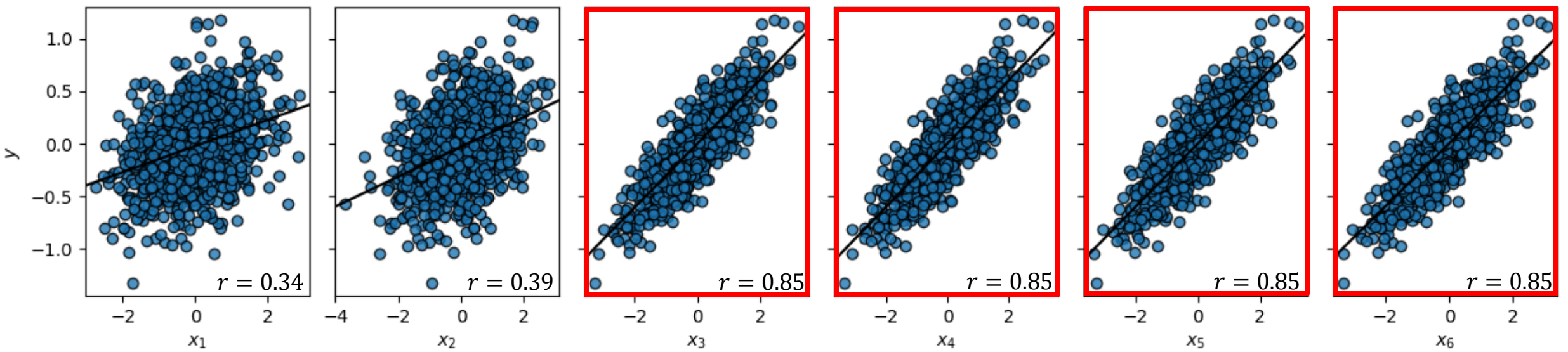
SIS or "Correlation Learning" Reduces Complexity

SIS+SO:

find $S = \{j_1, \dots, j_s\}$ where $|\tilde{\mathbf{x}}_j^T \mathbf{y}| \geq |\tilde{\mathbf{x}}_{j+1}^T \mathbf{y}|$ for $1 \leq j < p$
 and apply SO to sub-matrix $\beta_n^{SO}(\mathbf{X}_S, \mathbf{y})$
 computationally **efficient** for **small s**: $O(np + T_{SO}(k, n, s))$
 but **inconsistent** if s too small



$y = 0.5x_1 + 0.5x_2$



[Fan, J., Lv, J. (2008) *Sure independence screening* J. R. Stat. Soc. Ser. B 70(5)]

SISSO is an Iterative Correlation Learning Procedure

SISSO:

set $\mathbf{r}_1 = \mathbf{y}$

for $l = 1, \dots, k$:

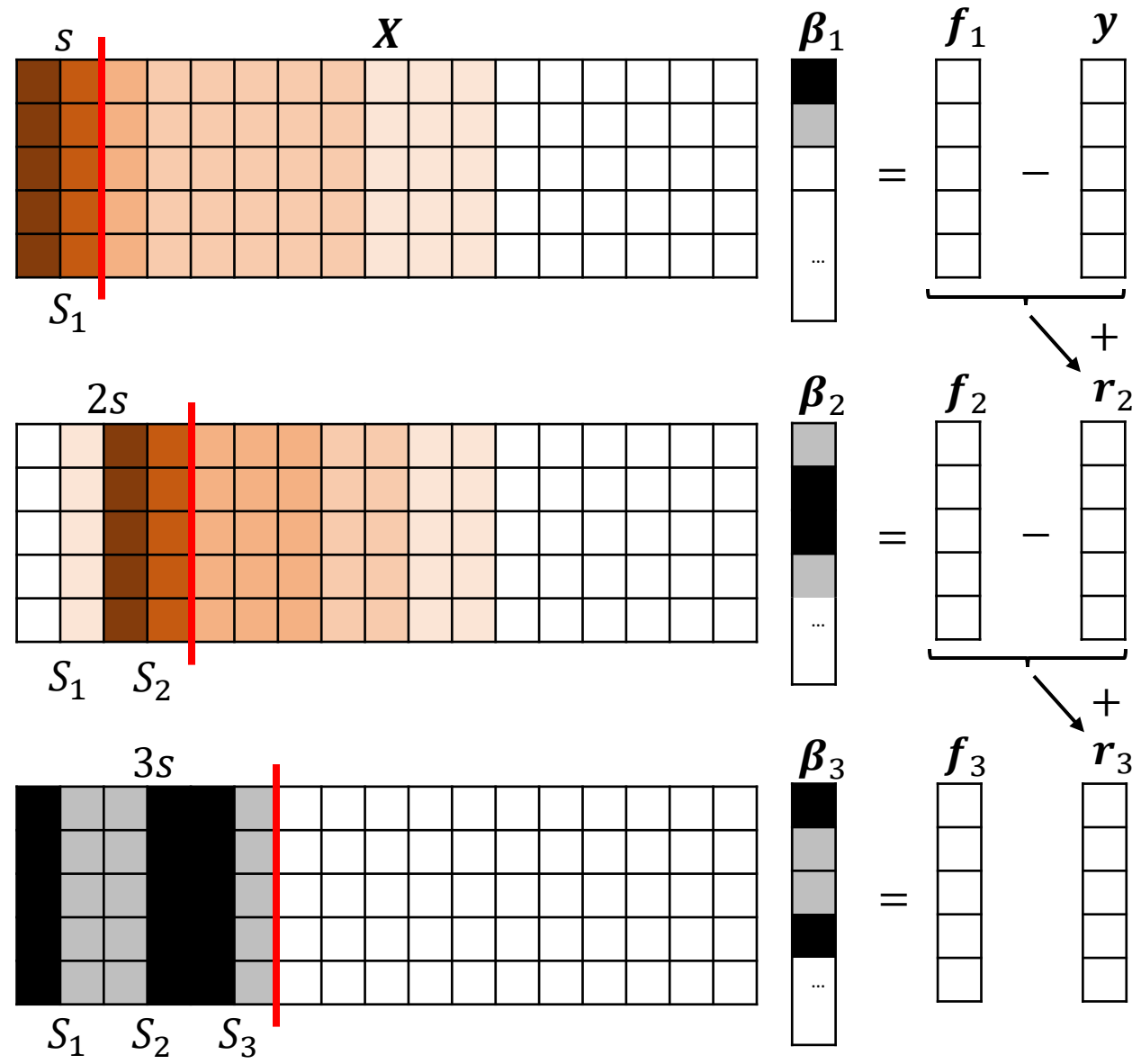
find $S_l = \{j_1, \dots, j_s\}$ s.t. $|\tilde{\mathbf{x}}_j^T \mathbf{r}_l| \geq |\tilde{\mathbf{x}}_{j+1}^T \mathbf{r}_l|$ for $1 \leq j < p$

set $\boldsymbol{\beta}_{l,n}^{\text{SISSO}} = \boldsymbol{\beta}_{l,n}^{\text{SO}}(\mathbf{X}_S, \mathbf{y})$ with $S = S_1 \cup \dots \cup S_l$

and $\mathbf{r}_{l+1} = \mathbf{y} - \mathbf{X}_S \boldsymbol{\beta}_{l,n}^{\text{SISSO}}$

Fundamental Questions:

1. What s is **computationally efficient**, i.e., what is s_{\max} st $T_{\text{ICL}}^{\text{SO}} \in O(knp + \sum_{l=1}^k T_{\text{SO}}(l, n, ls_{\max})) \leq c_0 + c_1 knp$?
2. What SO is **consistent** / performs best when choosing **optimal** $s \leq s_{\max}$?
3. Can performance be retained when choosing s **data-driven**?

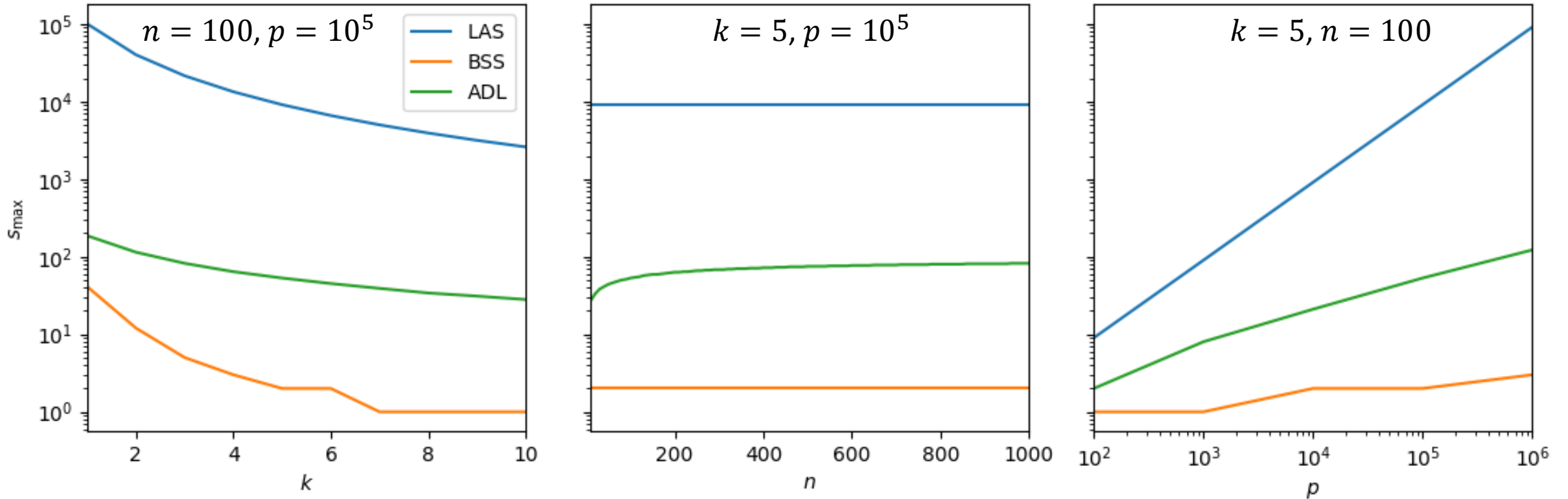


[Barut et al. (2016) *Conditional sure independence screening* JASA 111(515)]

[Fan, J., Lv, J. (2008) *Sure independence screening* J. R. Stat. Soc. Ser. B 70(5)]

[Ouyang et al. (2018) *SISSO: A compressed-sensing method for low-dimensional descriptors* Phys. Rev. Mater. 2(8)]

Computationally Feasible Pool Increment Values



General definition:

$$s_{\max}(k, n, p) = \max\{s \in \mathbb{N} : T_{\text{ICL}}(k, n, p, s) \leq c_0 + c_1 knp\}$$

Lasso:

$$s_{\max}^{\text{LAS}} \in \Theta(p/k^2)$$

Best-subset-search:

$$s_{\max}^{\text{BSS}} \in \Theta(\sqrt[k]{p})$$

Adaptive Lasso:

$$s_{\max}^{\text{BSS}} \in O\left(\min\left(\sqrt{p}, \sqrt[3]{np}\right)/k\right) \cap \Omega\left(\sqrt[3]{p/k^2}\right)$$

Evaluation over Wide Range of Functions

Ten correlated **normal primary inputs**

$$\mathbf{z} \sim N_{10}(\mathbf{0}, \mathbf{C}), C_{i,j} = 0.8^{|i-j|}$$

Degree $d = 1, 2, \dots, 7$ **multinomial feature maps**

$$\Phi_d = \{\boldsymbol{\varphi} \in \mathbb{N}^{10}: \|\boldsymbol{\varphi}\|_1 \leq d\}$$

$$x_{\boldsymbol{\varphi}} = \mathbf{z}^{\boldsymbol{\varphi}} = z_1^{\varphi_1} z_2^{\varphi_2} \dots z_{10}^{\varphi_{10}}$$

$$\mathbf{x} = (z_1^d, z_1^{d-1} z_2, z_1^{d-2} z_2 z_3, \dots, z_{10}^2, z_1, \dots, z_9, z_{10})$$

Random **sparse polynomials**

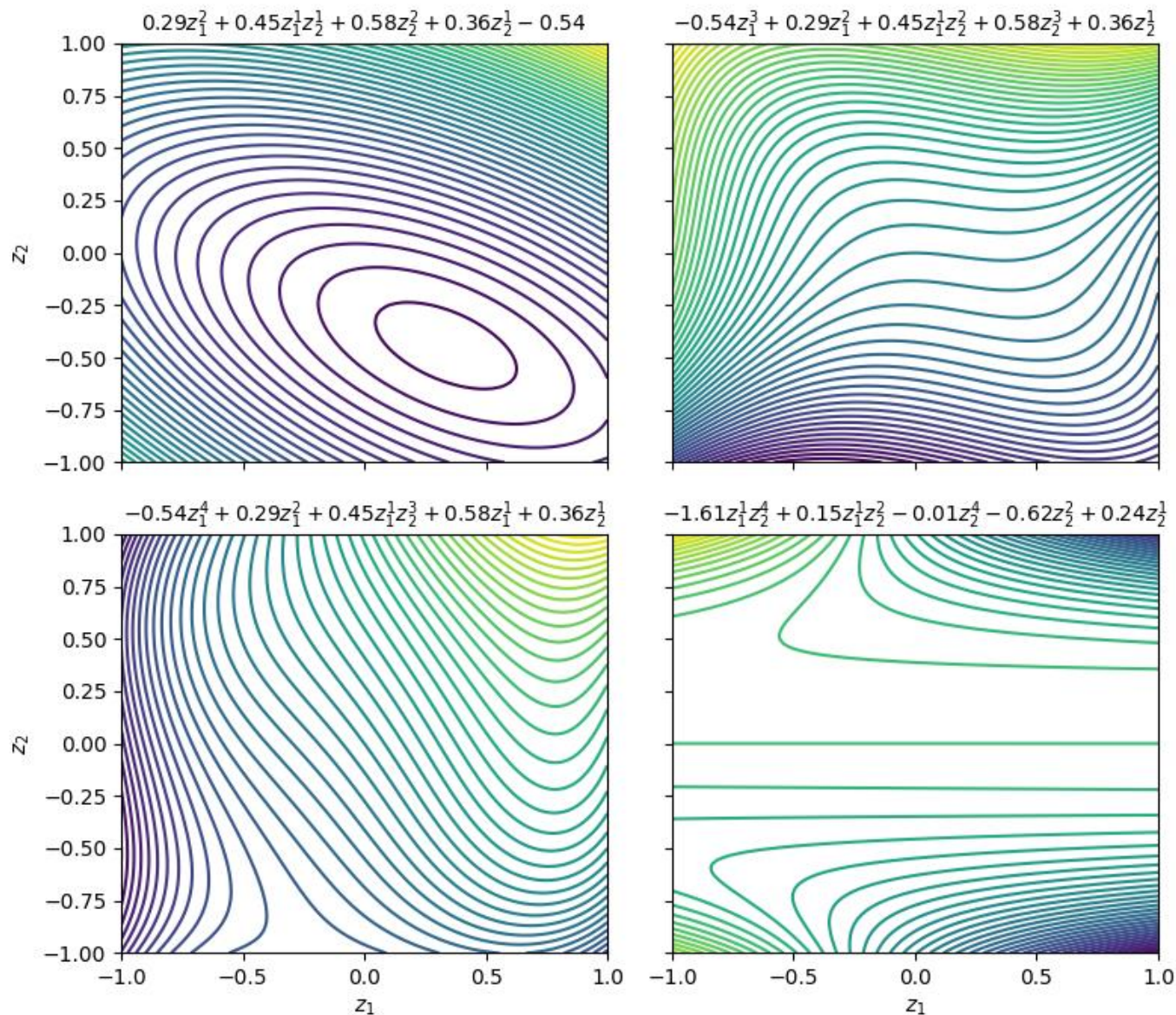
$$R = \{\boldsymbol{\varphi} \in \Phi: \varphi_6 = \dots = \varphi_{10} = 0\}$$

$$I^* \sim \text{Unif}(\{I \subseteq R: \#I = 5\})$$

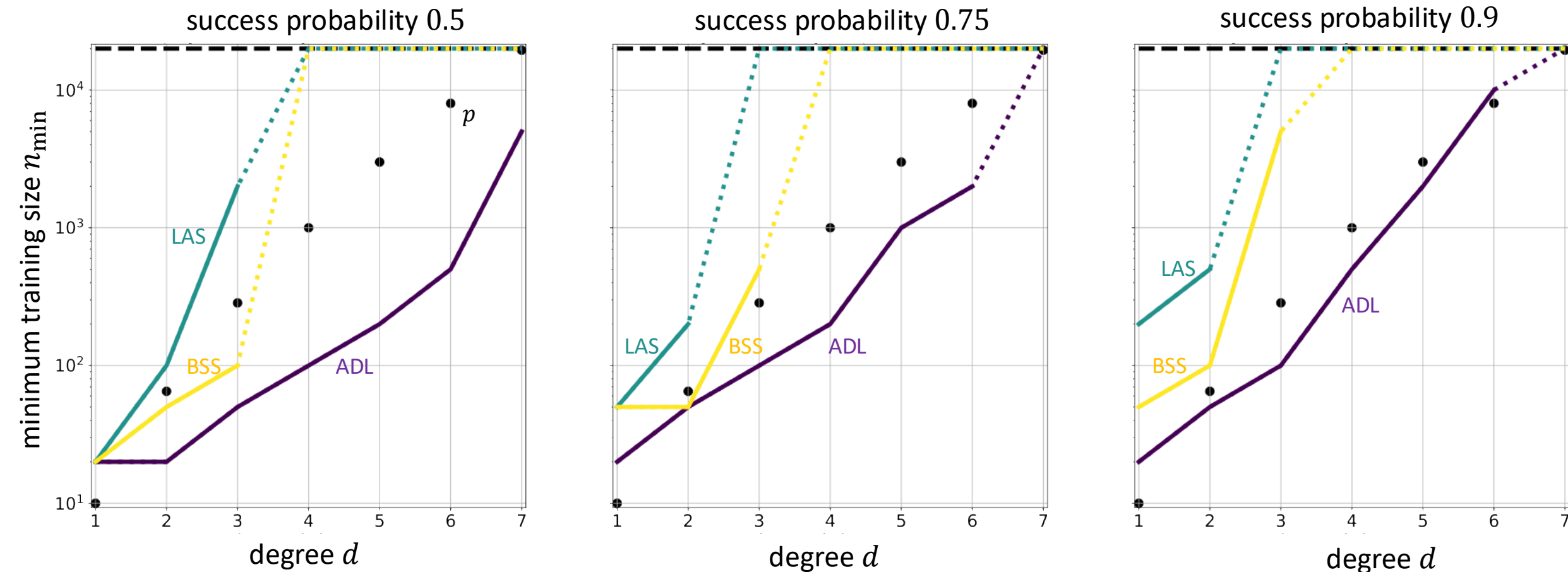
$$\beta_j^* \sim N(0, \sigma_j^{-1}) \text{ for } j \in I^* \text{ and } \beta_j^* = 0 \text{ for } j \notin I^*$$

Ten polynomials per degree

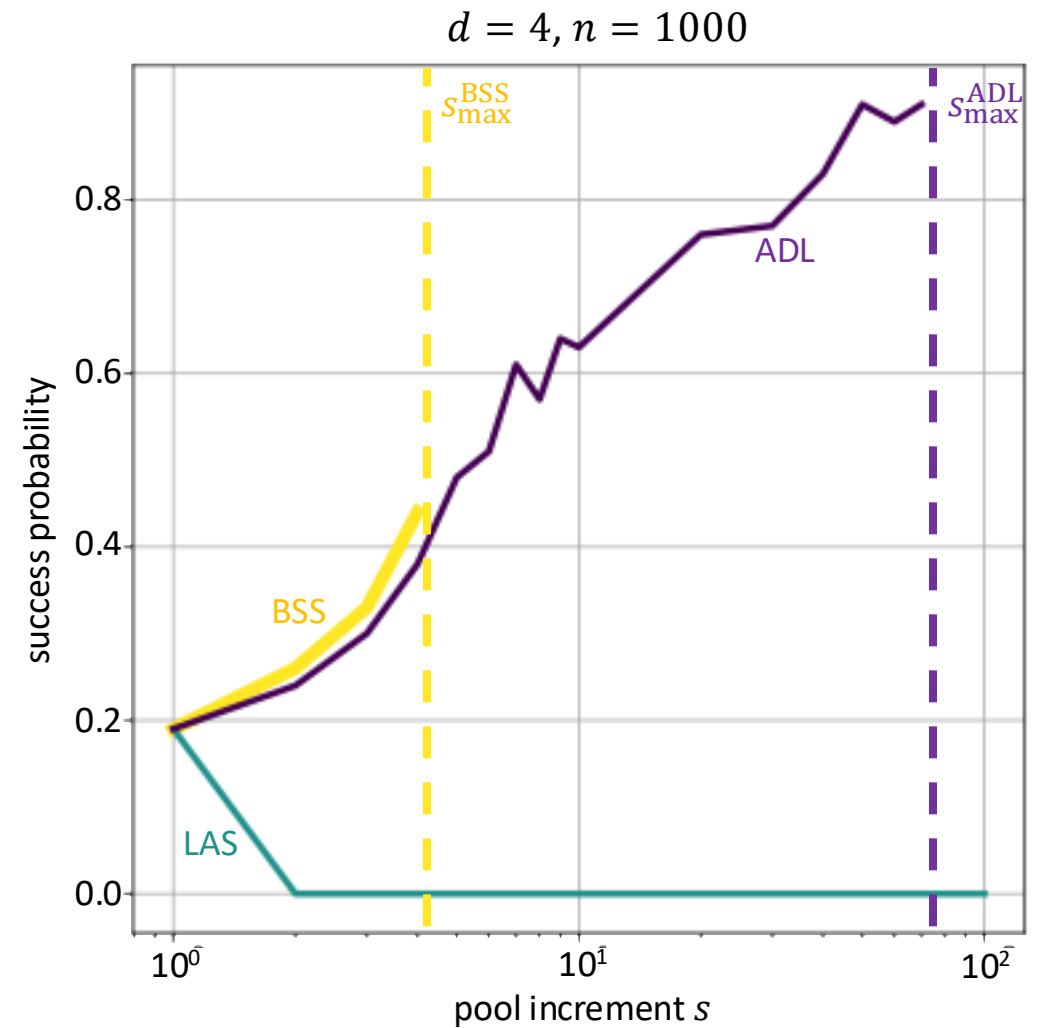
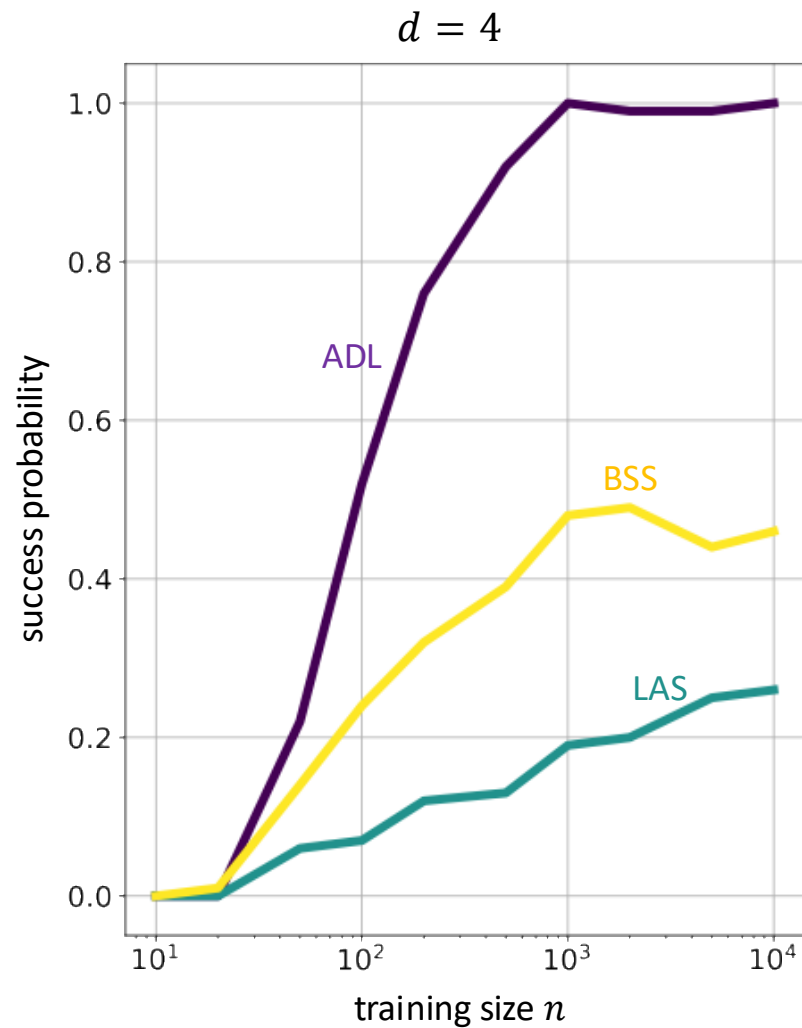
Ten datasets per polynomial



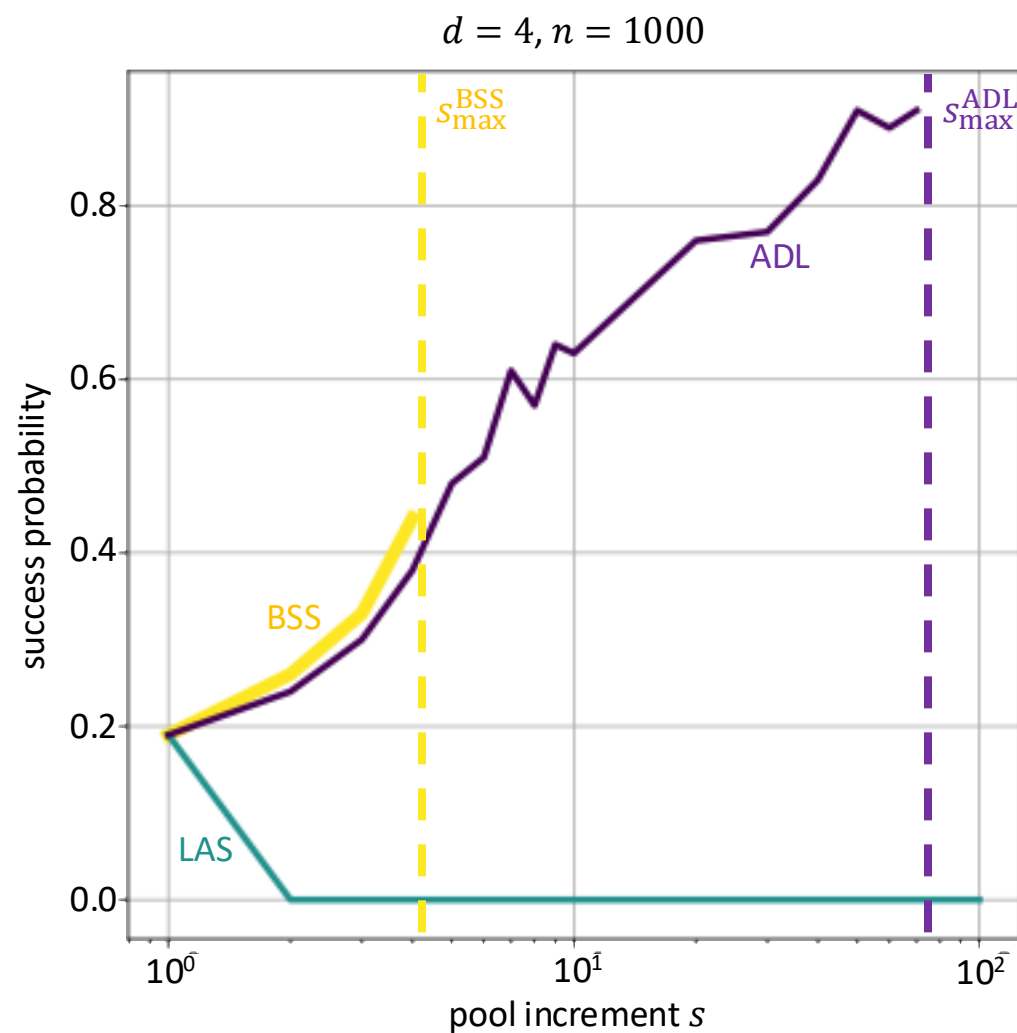
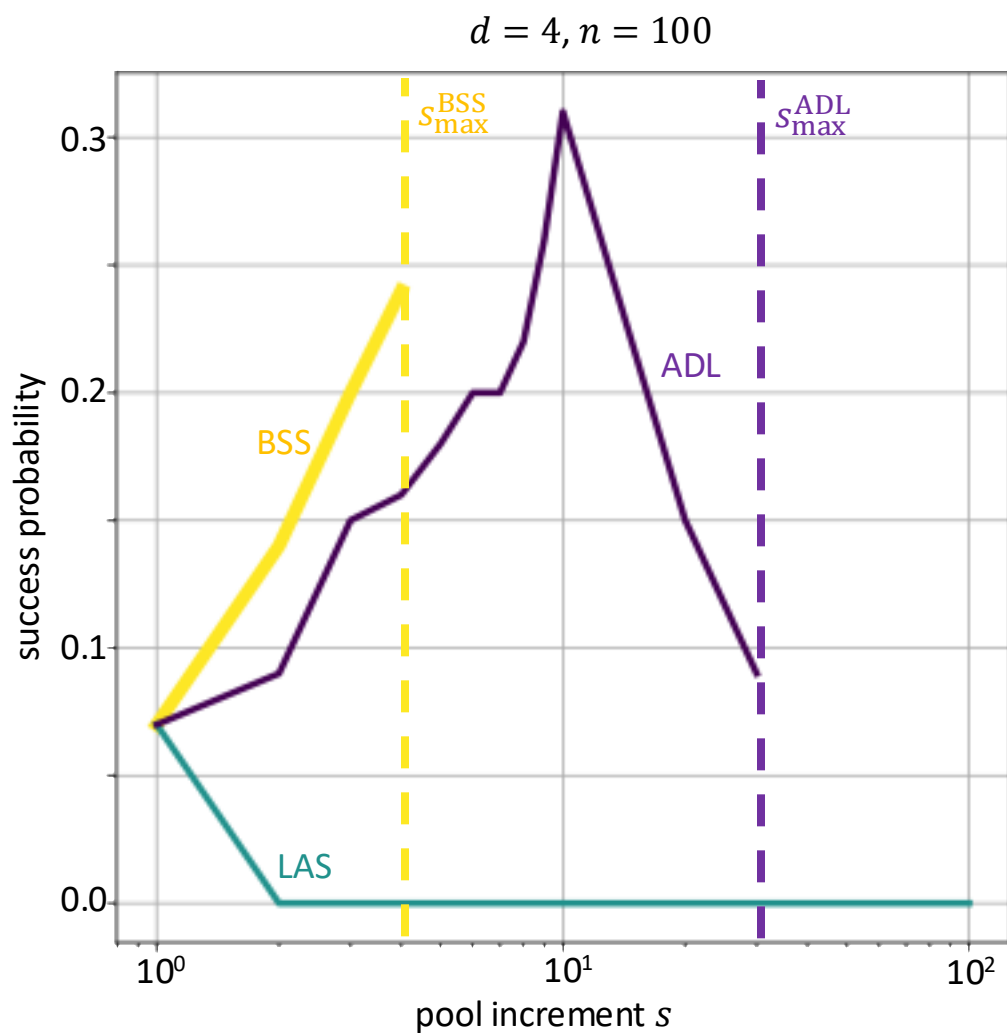
Adaptive Lasso Best-performing Sparsifying Operator 13



Advantage due Larger Range of Available s values



Maximum Pool Increment is not Always Optimal



Advantage Retained with Data-driven Selection

In practice: s_* **unknown** and s needs to be selected based on fixed rule or data, e.g., via **cross validation**:

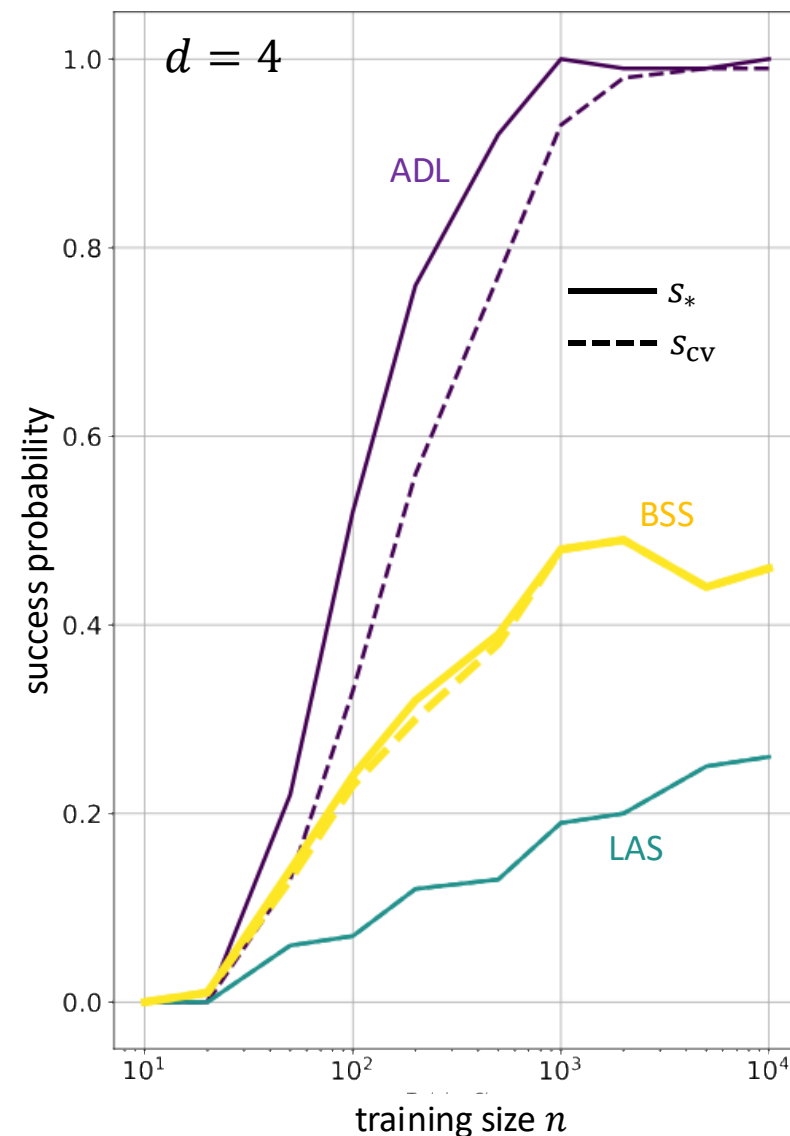
- $s_{cv} = \operatorname{argmin}\{\sum_{l=1}^{10} \|\mathbf{X}_l \boldsymbol{\beta}_l - \mathbf{y}_l\|^2 : 1 \leq s \leq s_{\max}\}$
- $\boldsymbol{\beta}_l = \boldsymbol{\beta}(\mathbf{X}_{\setminus l}, \mathbf{y}_{\setminus l}, s)$

Note: selection problem **hardest for adaptive Lasso**

- **BSS:** only few feasible s and $s = s_{\max}$ tends to work well
- **Lasso:** generally want very small s (1 or 2), i.e., slightly relaxed matching pursuit works better than Lasso
- **Adaptive Lasso:** relatively wide range available and need to trade off selection of relevant versus irrelevant variables

Result

- While data-driven selection reduces adaptive Lasso performance, marked **advantage retained** over BSS
- ...at least for degree 4 polynomials (limit due to 10x comp. cost)



Conclusion

Summary

- Investigate **identification consistency** and convergence rates of SISSO methods under explicit **computational constraint**
- **Adaptive Lasso** appears to be attractive SO, combining consistency with relative computational efficiency
- Indeed, **outperforms BSS and Lasso** in wide range of practical problems and retained when using **cross validation** to choose pool increment

Future

- **Theoretical bounds** for SISSO success probability
- Translation to **materials properties** modelling
- **Sparse regression estimators** with computational cost between ADL and BSS, e.g., SCAD, Dantzig Selector, iterative thresholding?

