

Quantifying Uncertainty and Stability Among Highly Correlated Predictors: A Subspace Perspective

Xiaozhu Zhang

Department of Statistics
University of Washington

Joint with



Armeen Taeb
University of Washington



Jacob Bien
University of Southern California

Setup

Suppose that we have a dataset containing (X, Y) , where

- $X \in \mathbb{R}^{n \times p}$ is a fixed design matrix for features
- Some features are **nearly linearly dependent**
- $Y \in \mathbb{R}^n$ is a vector for response variable

Model selection task

Which model $S \subseteq \{1, \dots, p\}$ can explain the response variable?

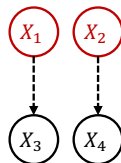
Table of Contents

- 1 Three challenges in highly-correlated variable selection
- 2 The subspace perspective
- 3 FSSS algorithm
- 4 Real data application

Three Challenges

Toy example:

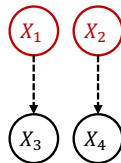
- $X_1 \approx X_3, \quad X_2 \approx X_4$
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$
- So the true support is: $S^* = \{1, 2\}$



Three Challenges

Toy example:

- $X_1 \approx X_3, \quad X_2 \approx X_4$
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$
- So the true support is: $S^* = \{1, 2\}$



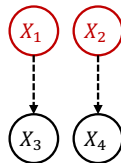
Challenge 1: Measuring Subspace Accuracy

How to define **True Positives (TP)** and **False Positives (FP)**?

Three Challenges

Toy example:

- $X_1 \approx X_3, \quad X_2 \approx X_4$
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$
- So the true support is: $S^* = \{1, 2\}$



Challenge 1: Measuring Subspace Accuracy

How to define **True Positives (TP)** and **False Positives (FP)**?

Suppose the selected set is: $\hat{S} = \{1, 4\}$

- Naively: 1 true positive, 1 false positive
- Since $X_4 \approx X_2$, is X_4 really “false”?

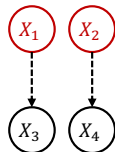
Goal: Redefine TP / FP so that

$$\text{TP}(S^*, \hat{S}) \approx 2, \quad \text{FP}(S^*, \hat{S}) \approx 0$$

Three challenges

Toy example

- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



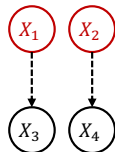
Challenge 2

How to quantify **stability** of selected features and sets?

Three challenges

Toy example

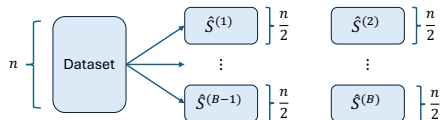
- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



Challenge 2

How to quantify **stability** of selected features and sets?

Stability selection (Meinshausen and Bühlmann, 2010):



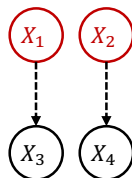
Stability of feature j :

- **selection proportion** $\pi(\{j\})$

Three Challenges

Toy example:

- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



Challenge 2

How to quantify **stability** of selected features and sets?

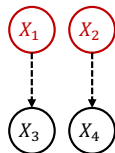
Stability selection (Meinshausen & Bühlmann, 2010):

- $\sim 50\%$ of subsamples select X_1
- $\sim 50\%$ of subsamples select X_3
- $\pi(\{1\}) \approx \pi(\{3\}) \approx 0.5 \Rightarrow$ **Not stable!**

Three challenges

Toy example

- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



Challenge 2

How to quantify **stability** of selected features and sets?

Stability selection (Meinshausen and Bühlmann, 2010):

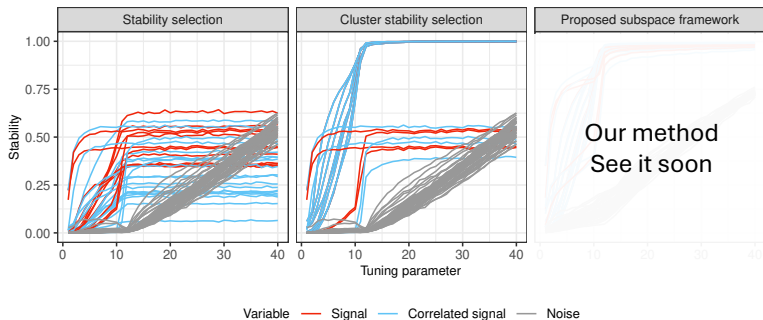
Use Lasso or ℓ_0 as the base procedure:

- $\sim 50\%$ subsamples choose X_1
- $\sim 50\%$ subsamples choose X_3
- $\pi(\{1\}) \approx \pi(\{3\}) \approx 0.5 \Rightarrow$ **NOT stable!**

Goal: re-define stability
s.t. $\pi(\{1\}) \approx \pi(\{3\}) \approx 1$

An experiment: Quantifying stability

A comprehensive synthetic dataset:



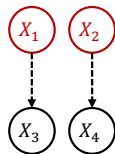
For stability selection (and its variant):

- Many signal and correlated signal features are **not stable**!

Three challenges

Toy example

- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



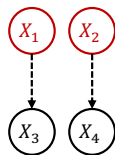
Challenge 3

How to **aggregate** features into a model?

Three challenges

Toy example

- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



Challenge 3

How to **aggregate** features into a model?

Suppose now all 4 variables are stable.

Can we have $\hat{S} = \{1, 2, 3, 4\}$?

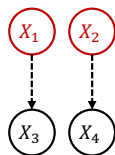
- No! Because X_1 and X_3 are redundant! (as well as X_2 and X_4)
- Instead, we want multiple models: $\{1, 2\}$, $\{1, 4\}$, $\{2, 3\}$, $\{3, 4\}$

Goal: an algorithm that outputs **multiple** models **without redundancy**

Beyond feature perturbation

Toy example ♠:

- X_1 and X_3 are highly correlated
 X_2 and X_4 are highly correlated
- $y = \beta_1^* X_1 + \beta_2^* X_2 + \epsilon$



We also need to deal with **more complicated** structures...

Toy example ♣:

- $X_3 = X_1 + X_2 + \delta$ (a perturbed **linear combination**)
- $y = X_1 - X_2 + \epsilon$

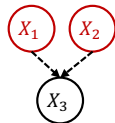


Table of Contents

- 1 Three challenges in highly-correlated variable selection
- 2 The subspace perspective
- 3 FSSS algorithm
- 4 Real data application

Our contribution: a subspace perspective

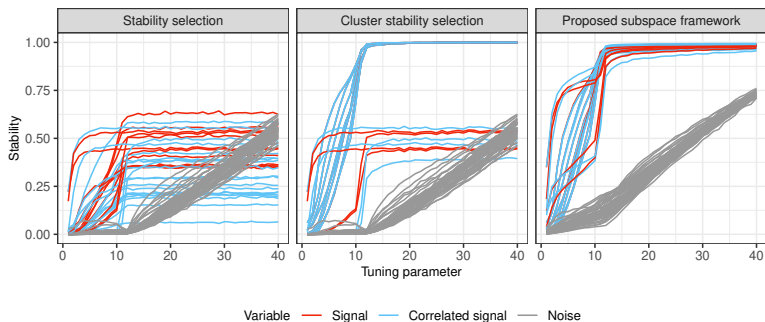
For the given fixed design matrix $X \in \mathbb{R}^{n \times p}$,

- **map** a selection set S onto its **column space** $\text{col}(X_S)$
- $\text{col}(X_S)$ is a subspace living in \mathbb{R}^n
- $\mathcal{P}_{X_S}(y)$ is the best linear prediction of y by S

Our contribution: a subspace perspective

For the given fixed design matrix $X \in \mathbb{R}^{n \times p}$,

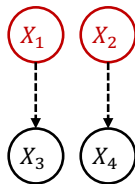
- **map** a selection set S onto its **column space** $\text{col}(X_S)$
- $\text{col}(X_S)$ is a subspace living in \mathbb{R}^n
- $\mathcal{P}_{X_S}(y)$ is the best linear prediction of y by S



Why Are Subspaces Useful?

In the toy example:

- Sets $\{1\}$ and $\{3\}$ share no features, yet $\text{col}(X_1) \approx \text{col}(X_3)$
- Sets $\{1, 2\}$ and $\{1, 4\}$ overlap in only one feature, yet $\text{col}(X_{1,2}) \approx \text{col}(X_{1,4})$

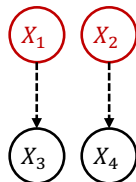


Example: $X_3 \approx X_1$, $X_4 \approx X_2$

Why Are Subspaces Useful?

In the toy example:

- Sets $\{1\}$ and $\{3\}$ share no features, yet $\text{col}(X_1) \approx \text{col}(X_3)$
- Sets $\{1, 2\}$ and $\{1, 4\}$ overlap in only one feature, yet $\text{col}(X_{1,2}) \approx \text{col}(X_{1,4})$



Example: $X_3 \approx X_1$, $X_4 \approx X_2$

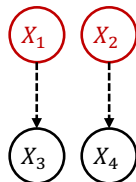
Key Idea

Subspace alignment captures similarity even when sets differ

Why Are Subspaces Useful?

In the toy example:

- Sets $\{1\}$ and $\{3\}$ share no features, yet $\text{col}(X_1) \approx \text{col}(X_3)$
- Sets $\{1, 2\}$ and $\{1, 4\}$ overlap in only one feature, yet $\text{col}(X_{1,2}) \approx \text{col}(X_{1,4})$



Example: $X_3 \approx X_1$, $X_4 \approx X_2$

Key Idea

Subspace alignment captures similarity
even when sets differ

Implications:

- Higher $\text{corr}(X_1, X_3) \Rightarrow$ stronger alignment: $\text{col}(X_1) \approx \text{col}(X_3)$
- Subspace alignment varies **smoothly** with correlation
- If features are **orthogonal**,
then alignment = exact set overlap

Subspace True positives and False Positives: Definitions

Definition (true positive, false positive)

Let S^* be the true set of features, and $\hat{S} \subseteq \{1, \dots, p\}$ be the estimated set. Then:

$$\begin{aligned}\text{TP}(\hat{S}, S^*) &:= \text{trace} \left(\mathcal{P}_{\text{col}(X_{\hat{S}})} \mathcal{P}_{\text{col}(X_{S^*})} \right), \\ \text{FPE}(\hat{S}, S^*) &:= |\hat{S}| - \text{TP}(\hat{S}, S^*)\end{aligned}$$

Subspace True positives and False Positives: Definitions

Definition (true positive, false positive)

Let S^* be the true set of features, and $\hat{S} \subseteq \{1, \dots, p\}$ be the estimated set. Then:

$$\begin{aligned}\text{TP}(\hat{S}, S^*) &:= \text{trace} \left(\mathcal{P}_{\text{col}(X_{\hat{S}})} \mathcal{P}_{\text{col}(X_{S^*})} \right), \\ \text{FPE}(\hat{S}, S^*) &:= |\hat{S}| - \text{TP}(\hat{S}, S^*)\end{aligned}$$

- TP measures how well the selected subspace aligns with the true subspace
- FPE counts the dimensions in \hat{S} that don't align with the true subspace
- Always: $\text{TP} \leq \min\{|\hat{S}|, |S^*|\}$

Subspace True positives and False Positives: Definitions

Definition (true positive, false positive)

Let S^* be the true set of features, and $\hat{S} \subseteq \{1, \dots, p\}$ be the estimated set. Then:

$$\begin{aligned}\text{TP}(\hat{S}, S^*) &:= \text{trace} \left(\mathcal{P}_{\text{col}(X_{\hat{S}})} \mathcal{P}_{\text{col}(X_{S^*})} \right), \\ \text{FPE}(\hat{S}, S^*) &:= |\hat{S}| - \text{TP}(\hat{S}, S^*)\end{aligned}$$

- TP measures how well the selected subspace aligns with the true subspace
- FPE counts the dimensions in \hat{S} that don't align with the true subspace
- Always: $\text{TP} \leq \min\{|\hat{S}|, |S^*|\}$

Geometric interpretation

$$\text{TP} = \sum_i (\cos \theta_i)^2 \quad \text{FPE} = |\hat{S}| - \sum_i (\cos \theta_i)^2$$

where θ_i are the **principal angles** between $\text{col}(X_{\hat{S}})$ and $\text{col}(X_{S^*})$.

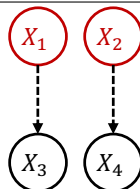
Subspace True Positives and False Positives: Example

$$\text{TP} = \sum_i (\cos \theta_i)^2 \quad \text{FPE} = |\hat{S}| - \sum_i (\cos \theta_i)^2$$

where θ_i are the **principal angles** between $\text{col}(X_{\hat{S}})$ and $\text{col}(X_{S^*})$.

- Let $S^* = \{1, 2\}$ and we estimate $\hat{S} = \{3, 4\}$
- In this example:

$$\theta_1 \approx 0^\circ, \quad \theta_2 \approx 0^\circ \quad \Rightarrow \quad \text{TP} \approx 2, \quad \text{FPE} \approx 0$$



Example: $X_3 \approx X_1, X_4 \approx X_2$

Subspace True Positives and False Positives: Example

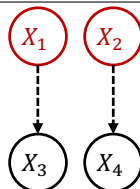
$$\text{TP} = \sum_i (\cos \theta_i)^2 \quad \text{FPE} = |\hat{S}| - \sum_i (\cos \theta_i)^2$$

where θ_i are the **principal angles** between $\text{col}(X_{\hat{S}})$ and $\text{col}(X_{S^*})$.

- Let $S^* = \{1, 2\}$ and we estimate $\hat{S} = \{3, 4\}$

- In this example:

$$\theta_1 \approx 0^\circ, \quad \theta_2 \approx 0^\circ \quad \Rightarrow \quad \text{TP} \approx 2, \quad \text{FPE} \approx 0$$



Example: $X_3 \approx X_1, X_4 \approx X_2$

More generally:

- principal angles vary smoothly
- if subspaces are orthogonal:

$$\text{TP} = |\hat{S} \cap S^*|, \quad \text{FPE} = |\hat{S} \setminus S^*| \quad \text{Reduces to classical notions}$$

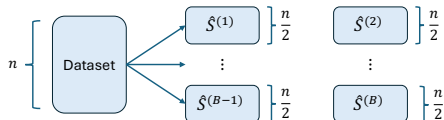
Subspace Stability (2nd challenge): Setup

Goal

Measure how consistently subspaces are selected across random subsamples.

Apply Lasso or ℓ_0 regression to:

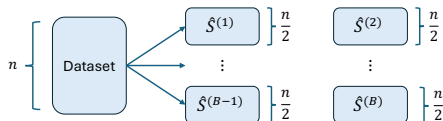
- B subsamples of size $\lfloor n/2 \rfloor$
- Each produces a set $\hat{S}^{(\ell)}$



Subspace Stability (2nd challenge): Setup

Apply Lasso or ℓ_0 regression to:

- B subsamples of size $\lfloor n/2 \rfloor$
- Each produces a set $\hat{S}^{(\ell)}$



Goal

Measure how consistently subspaces are selected across random subsamples.

Definition (stability)

Given B subsamples:

$$\mathcal{P}_{\text{avg}} := \frac{1}{B} \sum_{\ell=1}^B \mathcal{P}_{\text{col}(X_{\hat{S}^{(\ell)}})}.$$

Then for any set S ,

$$\begin{aligned} \pi(S) &:= \sigma_{|S|}(\mathcal{P}_{\text{col}(X_S)} \mathcal{P}_{\text{avg}} \mathcal{P}_{\text{col}(X_S)}) \\ &= \min_{\substack{z \in \text{col}(X_S) \\ \|z\|_2=1}} \frac{1}{B} \sum_{\ell=1}^B \|\mathcal{P}_{\text{col}(X_{\hat{S}^{(\ell)}})}(z)\|_2^2 \end{aligned}$$

What Does Stability $\pi(S)$ Measure?

Stability of set S with respect to estimates $\{\hat{S}^{(\ell)}\}_{\ell=1}^B$

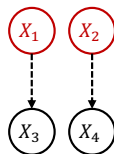
$$\pi(S) := \sigma_{|S|}(\mathcal{P}_{\text{col}(X_S)} \mathcal{P}_{\text{avg}} \mathcal{P}_{\text{col}(X_S)}).$$

with $\mathcal{P}_{\text{avg}} := \frac{1}{B} \sum_{\ell=1}^B \mathcal{P}_{\text{col}(X_{\hat{S}^{(\ell)}})}$

- $\pi(S) \in [0, 1]$
- $\pi(S)$ captures the **lowest alignment** of any direction in $\text{col}(X_S)$ with \mathcal{P}_{avg}
- $\pi(S)$ is the smallest squared cosine between any direction in $\text{col}(X_S)$ and the average of the subsample subspaces.
- High $\pi(S)$ means that $\text{col}(X_S)$ is reliably recovered across subsamples

Back to Toy Example: $\pi(\{1\}) \approx \pi(\{3\}) \approx 1$

- $\sim 50\%$ of subsamples select X_1
- $\sim 50\%$ select X_3
- But $X_3 = X_1 + \delta$ — nearly the same direction
- So all subsamples nearly capture the X_1 and X_3 direction!
- Therefore: $\pi(\{1\}) \approx \pi(\{3\}) \approx 1$

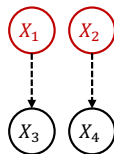


$$X_3 = X_1 + \delta \text{ } (\delta \text{ small})$$

Both signal features and highly correlated non-signal features have high stability π !

Back to Toy Example: $\pi(\{1, 3\}) \approx 0$

- $\sim 50\%$ of subsamples select X_1
- $\sim 50\%$ select X_3
- $\text{col}(X_{1,3})$ includes directions X_1 and δ
- But no subsample selects both X_1 and X_3 together
- So \mathcal{P}_{avg} misses the δ direction
- Therefore: $\pi(\{1, 3\}) \approx 0$



$$X_3 = X_1 + \delta \text{ } (\delta \text{ small})$$

Redundant features in S lead to small stability $\pi(S)$

Table of Contents

- 1 Three challenges in highly-correlated variable selection
- 2 The subspace perspective
- 3 FSSS algorithm
- 4 Real data application

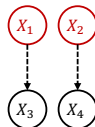
FSSS algorithm (3rd challenge)

We propose an algorithm “Features Subspace Stability Selection” (FSSS)

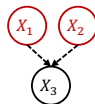
FSSS algorithm

- Input: Matrix \mathcal{P}_{avg} , and stability threshold $\alpha \in (1/2, 1)$
 - Sequentially add variables:
 - Start from the null model \emptyset
 - Add X_j to the current model \mathcal{S} if $\pi(\mathcal{S} \cup \{j\}) \geq \alpha$
 - Stop until no new features can be added
 - Output: one selection set $\hat{\mathcal{S}}$
-
- Any output $\hat{\mathcal{S}}$ is at least α -**stable**: $\pi(\hat{\mathcal{S}}) \geq \alpha$.
 - Each run may given different selection sets:
 - For example ♠, can return $\{1, 2\}$, $\{1, 4\}$, $\{2, 3\}$, $\{3, 4\}$
 - For example ♣, can return $\{1, 3\}$, $\{1, 2\}$, $\{2, 3\}$

Toy example ♠



Toy example ♣



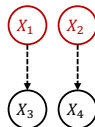
FSSS algorithm (3rd challenge)

We propose an algorithm “Features Subspace Stability Selection” (FSSS)

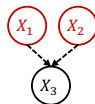
FSSS algorithm

- Input: Matrix \mathcal{P}_{avg} , and stability threshold $\alpha \in (1/2, 1)$
 - Sequentially add variables:
 - Start from the null model \emptyset
 - Add X_j to the current model \mathcal{S} if $\pi(\mathcal{S} \cup \{j\}) \geq \alpha$
 - Stop until no new features can be added
 - Output: one selection set $\hat{\mathcal{S}}$
-
- Any output $\hat{\mathcal{S}}$ is at least α -**stable**: $\pi(\hat{\mathcal{S}}) \geq \alpha$.
 - No redundant features can be included:
 - For example ♠ (as well as example ♣), $\{1, 2, 3\}$ can not be returned since $\pi(1, 2, 3) \approx 0$!

Toy example ♠



Toy example ♣

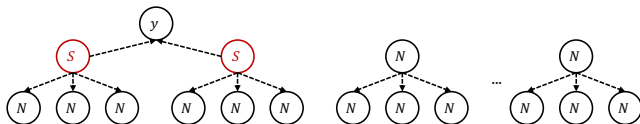


FSSS algorithm

Now that the models returned by FSSS are stable...
How accurate (FPE) can they be?

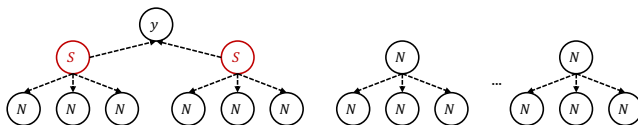
FSSS algorithm

Now that the models returned by FSSS are stable...
How accurate (FPE) can they be?



FSSS algorithm

Now that the models returned by FSSS are stable...
How accurate (FPE) can they be?



Theorem (False positive error control)

Under the clustering setup (as above), for any \hat{S} returned by FSSS with $|\hat{S}| \leq s_0$,

$$\mathbb{E} [\text{FPE}(\hat{S}, S^*)] \leq \frac{p(\gamma + b_n)^2}{2\alpha - 1} + a_n.$$

- The term γ is a **quality term** of the base procedure, with $\gamma \approx s_0/p$.
- a_n and b_n are **slackness** terms from dealing with perturbation and decomposing projection matrix.

FSSS algorithm

Now that the models returned by FSSS are stable...
How accurate (FPE) can they be?

Theorem (False positive error control)

Under the clustering setup (as above), for any \hat{S} returned by FSSS with $|\hat{S}| \leq s_0$,

$$\mathbb{E} [\text{FPE}(\hat{S}, S^*)] \leq \frac{p(\gamma + b_n)^2}{2\alpha - 1} + a_n.$$

- The term γ is a **quality term** of the base procedure, with $\gamma \approx s_0/p$.
 - a_n and b_n are **slackness** terms from dealing with perturbation and decomposing projection matrix.
-
- Stability-selection-type upper bound
 - orthogonal features \Rightarrow terms a_n and b_n **vanish** \Rightarrow UBD becomes $\frac{s_0^2}{p(1-2\alpha)}$

Experiments

A synthetic dataset including
...multiple **cluster** blocks, **parent-children** blocks, and **independent** features.

Base procedure: L0

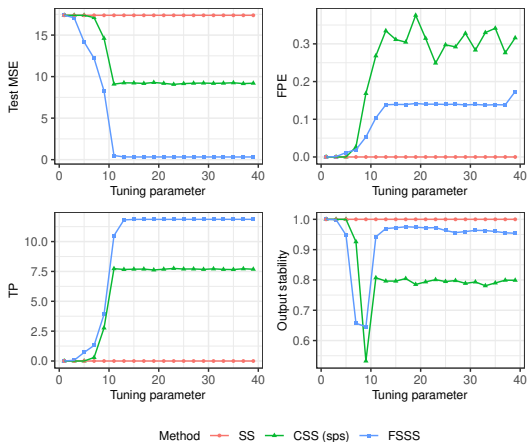


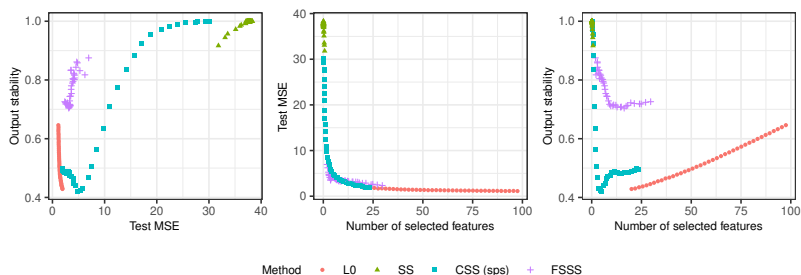
Table of Contents

- 1 Three challenges in highly-correlated variable selection
- 2 The subspace perspective
- 3 FSSS algorithm
- 4 Real data application

Gene expression in breast cancer

A gene expression dataset with $n = 189$, and $p = 1111$

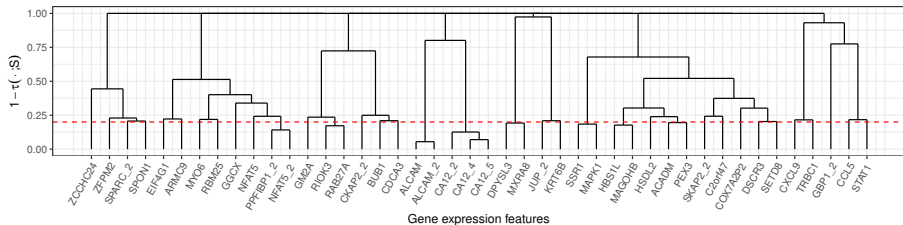
• Step 1: Performance comparison



• Step 2: Repeat FSSS algorithm to obtain 45 stable models

Gene expression in breast cancer

Among the 45 stable models:



Even more:

suppose that a domain expert give me some interesting features sets

We can do:

- rank those sets, or features, based on stability
- calibrate the commonality among those sets
- calibrate the substitutability among those sets

For model selection among highly-correlated predictors:

- Proposed a subspace framework:
 - Re-defined **true positive** & **false positive error**
 - Re-defined **stability**
 - Designed an algorithm **FSSS** that outputs stable models
- Future work:
 - More efficient FSSS algorithm
 - Extension to the non-linear setup: generalized additive models

Thanks for listening!

Check out the **full article**:

<https://arxiv.org/abs/2505.06760>

Check out the **package**:

<https://github.com/Xiaozhu-Zhang1998/substab>