# Comparison of Different Feaure Screening Methods with Application to BBS Associated Genes

Xiufang Wang

Institute of Statistics and Big Data
Renmin University of China

April 26, 2022

### Presentation Overview

- 1 Introduction
- 2 Feature Screening Marginal utility measures for screening Comparison of screening procedures
- 3 Results
- 4 Discussion and Future Work
- 6 Referencing

### Introduction

- 1 Data background: Bardet-Biedl syndrome (BBS) is a rare, inherited condition that can affect most organs in the body. So far, mutations in 21 genes have been identified as causing up to 80% of BBS cases. In the previous work, many researchers aim to find out additive genes having high association with the known causative genes based on the pairwise correlations. In this way, the corresponding gene therapy could be developed for the disease treatment.
- 2 Data source: from an eQTL experiment in rat eye loaded in <a href="https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE5680">https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE5680</a>.
- 3 Data description:
  - sample size: 120 twelve-week old male rats
  - dimensionality: 31099 different gene probe sets

### Introduction

### Our target:

- compare the screening performance with different marginal measures in terms of selecting related genes of TRIM 32( or BBS11)
- combine each screening method with group lasso algorithm and apply it to identify the potential genes having high association with BBS14, which causes up to 6% BBS disease.

Feature screening procedures can be roughly divided into two categories. The first is model-based.

- Sure independence screening (SIS):
  - based on Pearson correlation coefficient
  - for linear models
  - $\mathbb{E}(X_k^2) < \infty$  for  $k = 1, 2, \dots, p$
  - sensitive to outlies or extreme values in observations
- 2 Robust rank correlation screening (RRCS):
  - ullet based on Kendall au correlation coefficient
  - for semiparametric models and single-index models

cons: model misspecification

The second is model-free.

- 1 Sure independence ranking and screening (SIRS)
  - $\Omega(y) = \mathbb{E}\{\mathbf{x}F(y \mid \mathbf{x})\} = \text{cov}\{\mathbf{x}, \mathbf{1}(Y < y)\}.$
  - population quantity:

$$\omega_k = \mathbb{E}\{\Omega_k^2(Y)\}, k = 1, \dots, p$$

where  $\Omega_k(y)$  is the k th element of  $\Omega(y)$ .

sample version:

$$\tilde{\omega}_k = \frac{1}{n} \sum_{j=1}^n \left\{ \frac{1}{n} \sum_{i=1}^n X_{ik} \mathbf{1} (Y_i < Y_j) \right\}^2, \quad k = 1, \dots, p$$

where  $X_{ik}$  denotes the k th element of  $\mathbf{x}_i$ .

• robust to the outlies in the response observations.



- 2 Distance correlation based SIS (DC-SIS)
  - conditions: X and Y have finite first moments.
  - based on distance correlation  $\mathcal{R}(X, Y)$ :

$$\mathcal{R}^{2}(X, Y) = \begin{cases} \frac{\mathcal{V}^{2}(X, Y)}{\sqrt{\mathcal{V}^{2}(X)\mathcal{V}^{2}(Y)}}, & \mathcal{V}^{2}(X)\mathcal{V}^{2}(Y) > 0\\ 0, & \mathcal{V}^{2}(X)\mathcal{V}^{2}(Y) = 0 \end{cases}$$

where

$$V^{2}(X, Y) = \|f_{X,Y}(t, s) - f_{X}(t)f_{Y}(s)\|^{2}$$

- $f_X(\cdot)$ ,  $f_Y(\cdot)$  and  $f_{X,Y}(\cdot,\cdot)$  are characteristic functions of X, Y and (X,Y).
- for multivariate responses; screening grouped variables

- 3 Martingale distance correlation based SIS (MDC-SIS)
  - conditions:  $\mathbb{E}(|Y|^2 + |X|^2) < \infty$
  - based on martingale distance correlation MDC(Y | X):

$$[-\mathbb{E}\{Y - \mathbb{E}(Y)\}\{\tilde{Y} - \mathbb{E}(Y)\}\|X - \tilde{X}\|]^{1/2}/\{\operatorname{var}(Y)\operatorname{dvar}(X)\}^{1/2}$$

• sample version:  $\widehat{\omega}_k = MDC_n(Y \mid X_k)^2$ 

- 4 Culmulative divergence based SIS (CD-SIS)
  - conditions: var(X) > 0 and  $0 < var(Y) < \infty$
  - based on culmulative divergence CD(Y | X):

$$\mathbb{E}\left[\mathsf{cov}^2\{\mathit{Y},\mathbf{1}(\mathit{X}<\widetilde{\mathit{X}})\mid\widetilde{\mathit{X}}\}\right]/\mathsf{var}(\mathit{Y})$$

sample version:

$$\widehat{\mathsf{CD}}(Y \mid X) \stackrel{\mathsf{def}}{=} n^{-3} \sum_{j=1}^{n} \left[ \sum_{i=1}^{n} \left( Y_{i} - \bar{Y} \right) \left\{ I(X_{i} < X_{j}) - F_{n}(X_{j}) \right\} \right]^{2} / \widehat{\mathsf{var}}(Y)$$

- Generalize measure of correlation based SIS (GMC-SIS)
  - conditions: var(X) > 0 and  $0 < var(Y) < \infty$
  - based on generalize measure of correlation GMC(Y | X):

$$1 - \frac{\mathbb{E}\{Y - \mathbb{E}(Y \mid X)\}^2}{\mathsf{var}(Y)}$$

slicing estimator:

$$\widehat{\mathsf{GMC}}(Y \mid X) = 1 - \Lambda_n / \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

where

$$\Lambda_n = \frac{1}{H} \sum_{h=1}^{H} \frac{1}{c-1} \sum_{j=1}^{c} (Y_{(h,j)} - \frac{1}{c} \sum_{j=1}^{c} Y_{(h,j)})^2$$

• tuning parameter: c or equivalently H; the larger c is, the greater the performance of the GMC-SIS is.

- **6** A new correlation based SIS  $(\xi_n$ -SIS)
  - first rearrange the data as  $(X_{(1)},Y_{(1)}),\ldots,(X_{(n)},Y_{(n)})$  such that  $X_{(1)}\leq\ldots\leq X_{(n)}.$
  - based on

$$\xi(X,Y) := \frac{\int \operatorname{var}\left(\mathbb{E}\left(1_{\{Y \ge t\}} \mid X\right)\right) d\mu(t)}{\int \operatorname{var}\left(1_{\{Y \ge t\}}\right) d\mu(t)}$$

sample version:

$$\xi_n(X, Y) := 1 - \frac{3\sum_{i=1}^{n-1} |r_{i+1} - r_i|}{n^2 - 1}$$

where  $r_i$  is the rank of  $Y_{(i)}$ .

# Comparison of screening procedures

$$\mathcal{F}_1 = \{1 \le i \le p : \beta_i \ne 0\}$$

$$\mathcal{F}_2 = \{k : F(y \mid \mathbf{X}) \text{ functionally depends on } X_k \text{ for some } y \in \text{supp}(Y)\}$$

$$\mathcal{F}_3 = \{k : \mathbb{E}(y \mid \mathbf{X}) \text{ functionally depends on } X_k \text{ for some } y \in \text{supp}(Y)\}$$

Methods	Model-free	Active set	Standardization	Range
SIS	F	$\mathcal{F}_1$	F	[-1,1]
RRCS	F	$\mathcal{F}_1$	F	[-1,1]
SIRS	T	$\mathcal{F}_2$	T	$[0,\infty]$
DC-SIS	T	$\mathcal{F}_2$	F	[0,1]
MDC-SIS	T	$\mathcal{F}_3$	F	[0,1]
CD-SIS	T	$\mathcal{F}_3$	F	[0,1/4]
GMC-SIS	T	$\mathcal{F}_3$	F	[0,1]
$\xi_n$ -SIS	Т	$\mathcal{F}_2$	F	[0,1]

Table: Comparison of screening procedures

# Comparison of screening procedures

Methods	Robust of X	Invariance	Tuning param
SIS	F	scale change of $X$ and $Y$	F
RRCS	T	monotone trans of $X$ and $Y$	F
SIRS	F	monotone trans of $Y$	F
DC-SIS	F	scale change of $X$ and $Y$	F
MDC-SIS	F	scale change of $X$ and $Y$	F
CD-SIS	T	monotone trans of $X$	F
GMC-SIS	Т	monotone trans of $X$	slice number
$\xi_n$ -SIS	Т	monotone trans of $X$ and $Y$	F

Table: Comparison of screening procedures

Nine genes have been selected as the most relevant ones of gene BBS11 by INIS-penGAM algorithm in Fan, Feng, and Song (2011).

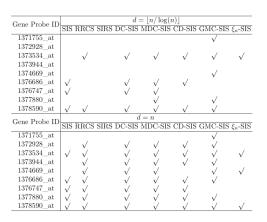


Figure: Results of gene screening using different methods

Combine each feature screening method with group lasso algorithm.

Methods	d	=500		d = n	$d = \lfloor$	$n/\log(n)$
Methods	Ave.S	Ave.MSPE	Ave.S	Ave.MSPE	Ave.S	Ave.MSPE
SIS	30.98	0.232	22.65	0.245	13.01	0.207
RRCS	31.5	0.231	22.48	0.256	13.42	0.261
SIRS	45.58	0.257	27.02	0.239	12.2	0.284
DC-SIS	45.52	0.244	13.92	0.109	5.77	0.168
MDC-SIS	27.37	0.246	25.66	0.302	8.62	0.309
CD-SIS	46.68	0.247	31.74	0.238	12.38	0.372
GMC-SIS(c=12)	29.41	0.239	24.15	0.248	13.09	0.268
GMC-SIS(c=20)	31.22	0.232	13.98	0.253	11.9	0.253
$\xi_n$ -SIS	33.08	0.241	22.58	0.256	12.36	0.230

Figure: Average model size(Ave.S) and average mean square prediction error(Ave.MSPE) with  $d = \lfloor n/\log(n) \rfloor$ , n and 500 out of 200 replications



Results of feature screening by different marginal dependence measures.

# Choose a marginal screening procedure:

#### Number of features to retain :

10

Option: number of observations within each slice for Sliced-GMC-SIS:

12

### Screening results

rank	gene	corr
1	X1383151_at	0.887262833947281
2	X1384172_at	0.885810091385537
3	X1373764_at	0.884794021898993
4	X1372197_at	0.882881006564438
5	X1392511_at	0.882193311529701
6	X1377829_at	0.879248051208977
7	X1373507_at	0.876312999268522
8	X1388896_at	0.874496094998033
9	X1389065_at	0.873854866228167
10	X1376067 at	0.871388010261845

#### Screening results

Screening results				
rank	gene	corr		
1	X1371823_at	0.710924369747899		
2	X1382045_at	0.696358543417367		
3	X1383151_at	0.69327731092437		
4	X1375887_at	0.678991596638655		
5	X1384172_at	0.677871148459384		
6	X1386952_a_at	0.673669467787115		
7	X1373764_at	0.669467787114846		
8.5	X1377829_at	0.667787114845938		
8.5	X1393214_at	0.667787114845938		
10	X1373507_at	0.66750700280112		



### Screening results

rank	gene	corr
1	X1383265_at	0.0673678453585641
2	X1371823_at	0.0664757480790294
3	X1375887_at	0.0662613788061506
4	X1383151_at	0.0656972027569693
5	X1375061_at	0.0654945847185489
6	X1390497_at	0.0652451861899022
7	X1398596_at	0.0647476544458326
8	X1373049_at	0.0645236813384127
9	X1389383_at	0.0642673728014739
10	X1373507_at	0.0641691138493283

Results of feeture screening by different marginal dependence measures.

Choose a marginal screening procedure:

CD-SIS

Number of features to retain:

10

Option: number of observations within each slice for Sliced-GMC-SIS:

### Screening results

<b>-</b>	9	
rank	gene	corr
1	X1382045_at	0.0660144768094535
2	X1383151_at	0.06458698447966
3	X1373507_at	0.0642416574412438
4	X1371823_at	0.0639770443167898
5	X1384172_at	0.0636988529147645
6	X1377829_at	0.0633905335198017
7	X1380174_at	0.063130685113085
8	X1375887_at	0.0630028045823302
9	X1389383_at	0.0628855641202612

X1373049 at 0.0625314863938655



### Screening results

rank	gene	corr
1	X1371823_at	0.570595180220849
2	X1375887_at	0.553302312660601
3	X1382045_at	0.542884922564067
4	X1389383_at	0.528925619834711
5	X1393214_at	0.52871727203278
6	X1384172_at	0.526008750607681
7	X1375061_at	0.51954996874783
8	X1374605_at	0.517049795124661
9	X1390497_at	0.507882491839711
10	X1390534_at	0.50767414403778

Results of feature screening by different marginal dependence measures.

Choose a marginal screening procedure:

Sliced-GMC-SIS

Number of features to retain:

10

Option: number of observations within each alice for Sliced-GMC-SIS:

### Screening results

•				
rank	gene	corr		
1	X1383151_at	0.730600283538146		
2	X1373764_at	0.727983708691608		
3	X1384172_at	0.717511452782615		
4	X1375887_at	0.709634814948908		
5	X1393214_at	0.708785020362474		
6	X1371823_at	0.702438259957917		
7	X1388766_at	0.698213542703655		
8	X1377829_at	0.693668919893442		
9	X1392511_at	0.692189576111425		
10	X1373049 at	0.689417781101181		

### Discussion and future work

- Combining DC-SIS and GMC-SIS with group lasso seem to have a superior performance among all methods when d=n or  $d=\lfloor n/\log(n)\rfloor$ . Superisingly, since the parameter c in GMC-SIS could be predetermined, we expect this method to have a better performance by increasing c properly. [Zhu, 1995] stated the effect of c on the convergence rate of the slicing estimator  $\widehat{\text{GMC}}$ .
- derive an iterative version for the GMC-based screening procedure and compare it with the counterpart of other methods.
- adopt other variable selection methods in the second stage.
- apply these procedures on more examples.

### References I



Fan, Jianqing and Lv, Jinchi (2008) Sure independence screening for ultrahigh dimensional feature space. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 70(5), 849–911.



Li, Gaorong and Peng, Heng and Zhang, Jun and Zhu, Lixing (2012) Robust rank correlation based screening.

The Annals of Statistics 40(3),1846–1877.



Fan, Jianqing and Feng, Yang and Song, Rui (2011) Nonparametric independence screening in sparse ultra-high-dimensional additive models.

Journal of the American Statistical Association 106(494),544-557.



Zhu, Li-Ping and Li, Lexin and Li, Runze and Zhu, Li-Xing (2011) Model-free feature screening for ultrahigh-dimensional data. Journal of the American Statistical Association 106(496),1464-1475.



Li, Runze and Zhong, Wei and Zhu, Liping (2012)
Feature screening via distance correlation learning.

Journal of the American Statistical Association 107(499),1129–1139.

### References II



Shao, Xiaofeng and Zhang, Jingsi (2014)

Martingale difference correlation and its use in high-dimensional variable screening. *Journal of the American Statistical Association* 109(507),1302–1318.



Zhou, Tingyou and Zhu, Liping and Xu, Chen and Li, Runze (2020) Model-free forward screening via cumulative divergence. *Journal of the American Statistical Association* 115(531),1393–1405.



Zheng, Shurong and Shi, Ning-Zhong and Zhang, Zhengjun (2012) Generalized measures of correlation for asymmetry, nonlinearity, and beyond. *Journal of the American Statistical Association* 107(499),1239–1252.



Chatterjee, Sourav(2021)

A new coefficient of correlation.

Journal of the American Statistical Association 116(536),2009–2022.



Zhu, Li-Xing and Ng, Kai W (1995)

Asymptotics of sliced inverse regression.

Statistica Sinica ,727–736.



# **Thanks**

Questions? Comments?