

# Sure Explained Variability and Independence Screening

Zhengjun Zhang

Department of Statistics  
University of Wisconsin  
Madison, WI 53706, USA

Co-authors: Min Chen, Yimin Lian, Zhao Chen

- 1 What've been done in independence screening
- 2 Sure Explained Variability: Definitions and Properties
- 3 Sure explained variability and independence screening (SEVIS)
- 4 Illustrative exmaples
- 5 Integrated genomic explained variability analyses of ovarian carcinoma
- 6 A nonparametric estimator for SEV

## Variable selection

- the least absolute shrinkage and selection operator (LASSO) in Tibshirani (1996), group LASSO in Yuan and Lin (2006), and the adaptive LASSO in Zou (2006);
- the smoothly clipped absolute deviation (SCAD) in Fan (1997), Fan and Li (2001);
- the elastic net in Zou and Hasties (2005);
- the Dantzig selector in Candes and Tao (2007);
- ....

## Independence screening

- Sure independence screening (SIS) based on Pearson's correlation: Fan and Lv (2008), Liu *et al.* (2013);
- Model based screening: Fan and Song (2010), Fan *et al.* (2011), Zhu *et al.* (2011), Song *et al.* (2012), Zhang *et al.* (2016);
- Model free screening: DC-SIS by Li *et al.* (2012), QaSIS by He *et al.* (2013), Q-SIS by Wu and Yin (2015);
- **Our approach:** based on explained variability and model free.

# Outline

- 1 What've been done in independence screening
- 2 Sure Explained Variability: Definitions and Properties**
- 3 Sure explained variability and independence screening (SEVIS)
- 4 Illustrative examples
- 5 Integrated genomic explained variability analyses of ovarian carcinoma
- 6 A nonparametric estimator for SEV

## A new look of a twice-told tale:

$$\text{var}(Y) = \text{var}(E(Y|X)) + E(\text{var}(Y|X)). \quad (2.1)$$

- $\text{Var}(E(Y|X))$  measures the spread of the conditional mean (center) of  $Y$  given  $X$
- $\text{Var}(E(Y|X))/\text{Var}(Y)$  can certainly be interpreted as the explained variance of  $Y$  by  $X$ .

$$\text{SEV}(Y|X) = \frac{\text{var}(E(Y|X))}{\text{var}(Y)} = 1 - \frac{E(\text{var}(Y|X))}{\text{var}(Y)} = 1 - \frac{E[\{Y - E(Y|X)\}^2]}{\text{var}(Y)}$$

Correlation ratios: Kendall and Stuart (1979), Doksum and Samarov (1995), Wang (2001), Zheng, Shi, and Zhang (2012)

## Properties of SEV

Suppose both  $E(X^2) < \infty$  and  $E(Y^2) < \infty$ . Then

(P.1) If  $Y = g(X)$ , *almost surely* (a.s.), then  $\text{SEV}(Y|X) = 1$ .

(P.2) If  $Y = g(X) + \varepsilon$ , then  
$$\text{SEV}(Y|X) = \text{var}(g(X)) / (\text{var}(g(X)) + \text{var}(\varepsilon)).$$

(P.3) If  $Y = ag(X) + b + \varepsilon$ , where  $g(x) = x$ ,  $a \neq 0$ , then  
$$\text{SEV}(Y|X) = \rho_{XY}^2.$$
 Here  $\rho_{XY}$  is Pearson's correlation coefficient.

(P.4) If  $\rho_{XY} \neq 0$ ,  $\text{SEV}(Y|X) \neq 0$ .

(P.5) If  $\rho_{XY} = -1$  or  $1$ ,  $\text{SEV}(Y|X) = 1$ .

(P.6)  $\text{SEV}(Y|X) \geq \rho_{XY}^2$ .

(P.7) If  $\text{SEV}(Y|X) = 0$ ,  $\rho_{XY} = 0$ .

## Example 1

*Suppose  $X$  is a standard normal random variable. Let  $Y = X^2$ . Then  $SEV(Y|X)=1$ , while  $\rho_{XY} = 0$ , and the distance correlation between  $Y$  and  $X$  is approximately 0.54 by a numerical calculation.*



# Outline

- 1 What've been done in independence screening
- 2 Sure Explained Variability: Definitions and Properties
- 3 Sure explained variability and independence screening (SEVIS)**
- 4 Illustrative exmaples
- 5 Integrated genomic explained variability analyses of ovarian carcinoma
- 6 A nonparametric estimator for SEV

## Notations

- $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ ,  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_p)$ ,  $\mathbf{X}_k = (X_{1k}, \dots, X_{nk})^T$
- $n$  is the sample size, the dimension  $p$  is much larger than the sample size  $n$ , i.e.  $n \ll p$ .
- $\sum_{i=1}^n X_{ik} = 0$  and  $\frac{1}{n-1} \sum_{i=1}^n X_{ik}^2 = 1$  for  $k = 1, \dots, p$ .

## SEVIS procedure

- The index set of active predictors:

$$\mathcal{M} = \{1 \leq k \leq p : \text{the predictor } X_k \text{ has contributed to the response } Y\}. \quad (3.1)$$

- Denoted by  $\omega_k = \text{SEV}(Y|X_k)$ ,  $k = 1, \dots, p$ . We define a new index set by

$$\mathcal{M}^* = \{1 \leq k \leq p : \omega_k > 0\}. \quad (3.2)$$

- Denoting  $\hat{\omega}_k$  as the estimator of  $\omega_k$ , we define a truncated active index set by

$$\hat{\mathcal{M}}^* = \{1 \leq k \leq p : \hat{\omega}_k \geq cn^{-\tau}\}. \quad (3.3)$$

- A practically workable active index set:

$$\hat{\mathcal{M}}_d^* = \{1 \leq k \leq p : \hat{\omega}_k \text{ is among the first } d \text{ largest estimates in } \hat{\mathcal{M}}^*\}. \quad (3.4)$$

## SEVIS algorithm

- 1. Normalize  $X_k$ ,  $k = 1, \dots, p$ , such that  $\sum_{i=1}^n X_{ik} = 0$  and  $(n-1)^{-1} \sum_{i=1}^n X_{ik}^2 = 1$ ;
- 2. Calculate  $\hat{\omega}_k = \text{SEV}(Y|X_k)$  between  $Y$  and candidate predictors  $X_k$ ,  $k = 1, \dots, p$ ;
- 3. Rank  $\hat{\omega}_k$  in a decreasing order, that is  $\hat{\omega}_{k_1} > \hat{\omega}_{k_2} > \dots > \hat{\omega}_{k_p}$ ;
- 4. Choose the first  $d[n/\log(n)]$ ,  $X_{k_1}, \dots, X_{k_{d[n/\log(n)]}}$ , as the active predictors.  $d$  is a given integer.

## Assumptions

- **(A1)**  $f(x, y)$  has a bounded support on  $(x, y) \in [s_x, S^x] \times [s_y, S^y]$ .  $f^X(x)$  is strictly positive, and is uniformly continuous on  $x \in [s_x, S^x]$ . The second derivative  $(f^X(x))''$  and  $(\phi^{Y|X}(x))''$  exist and are uniformly bounded on  $x \in [s_x, S^x]$ .
- **(A2)** The symmetric kernel function  $K(z)$  has a bounded support, and satisfies

$$\sup_{-\infty < z < \infty} |K(z)| < \infty, \quad \lim_{z \rightarrow \infty} |zK(z)| = 0, \quad \int z^2 |K(z)| dz < \infty.$$

The bandwidth  $h$  satisfies  $nh^3 \rightarrow \infty$  and  $nh^4 \rightarrow 0$ , as  $n \rightarrow \infty$ .

- **(A3)**  $\min_{k \in \mathcal{M}^*} \omega_k \geq 2cn^{-\tau}$ , where constant  $c > 0$  and  $0 \leq \tau < \frac{1}{6}$ .
- **(A4)** For given constants  $c > 0$ ,  $0 \leq \tau < \frac{1}{6}$ , and  $\log(p) = o(n^{1/3-2\tau})$ ,  $\liminf_{p \rightarrow \infty} (\min_{k \in \mathcal{M}^*} \omega_k - \max_{k \notin \mathcal{M}^*} \omega_k) > 2cn^{-\tau}$ .

## Sure Screening Property

Under Assumptions (A1),(A2), it follows that

$$\mathbb{P}\left(\max_{1 \leq k \leq p} |\hat{\omega}_k - \omega_k| \geq cn^{-\tau}\right) \leq O(p \exp\{-c_1 n^{1-2\tau} h^2\}). \quad (3.5)$$

Furthermore, if Assumption (A3) holds, we have that

$$\mathbb{P}\left(\mathcal{M}^* \subseteq \widehat{\mathcal{M}}^*\right) \geq 1 - O(s_n \exp\{-c_1 n^{1-2\tau} h^2\}), \quad (3.6)$$

where  $s_n$  is the cardinality of  $\mathcal{M}^*$ .

## Ranking Consistency property

Under Assumptions (A1), (A2), and (A4), we have that

$$\liminf_{n \rightarrow \infty} \{ \min_{k \in \mathcal{M}^*} \hat{\omega}_k - \max_{k \notin \mathcal{M}^*} \hat{\omega}_k \} \geq 0, \quad \text{a.s.} \quad (3.7)$$

# Outline

- 1 What've been done in independence screening
- 2 Sure Explained Variability: Definitions and Properties
- 3 Sure explained variability and independence screening (SEVIS)
- 4 Illustrative exmaples**
- 5 Integrated genomic explained variability analyses of ovarian carcinoma
- 6 A nonparametric estimator for SEV



## Settings and Example 1

- Existing methods to be compared: SIRS by Zhu et al. (2011), DC-SIS by Li et al. (2012) and Q-SIS with  $\alpha = 0.75$  by Wu and Yin (2015), i.e. **Model free**
- $d_1 = \lceil n/\log(n) \rceil$ ,  $d_2 = 2\lceil n/\log(n) \rceil$ , and  $d_3 = 3\lceil n/\log(n) \rceil$

$$(1.a) \quad Y = c_1 X_1 + c_1 X_2 + c_1 X_{12} + c_1 X_{22} + \varepsilon,$$

$$(1.b) \quad Y = c_1 X_1 + c_2 X_2 + c_3 \mathbf{I}(X_{12} < 0) + c_4 X_{22} + \varepsilon,$$

$$(1.c) \quad Y = c_1 X_1 X_2 + c_3 \mathbf{I}(X_{12} < 0) + c_4 X_{22} + \varepsilon,$$

$$(1.d) \quad Y = c_1 X_1 + c_2 X_2 + c_3 \mathbf{I}(X_{12} < 0) + c_5 \exp(|X_{22}|)\varepsilon.$$

The predictors  $X = (X_1, \dots, X_p)$  are generated from a normal distribution with mean zero and covariance matrix  $\Sigma = (\sigma_{ij})_{p \times p}$ , where  $\sigma_{ij} = \sigma^{|i-j|}$ . The random error  $\varepsilon$  follows a standard normal distribution. We set the vector  $(c_1, c_2, c_3, c_4, c_5) = (5, 2, 7, 5, 2)$ . The sample size  $n = 200$ , the dimension  $p = 2000, 5000$ , and  $\sigma = 0.5, 0.8$ .

**Table 1 :** The Proportions of  $\mathcal{P}_s$  and  $\mathcal{P}_a$ ,  $p = 2000$   $\sigma = 0.5$

Model	Size	SIRS	DC-SIS	Q-SIS	SEVIS				
		$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$				$\mathcal{P}_a$
		<i>ALL</i>	<i>ALL</i>	<i>ALL</i>	$X_1$	$X_2$	$X_{12}$	$X_{22}$	<i>ALL</i>
(1.a)	$d_1$	<b>1.00</b>	<b>1.00</b>	<b>0.98</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_2$	<b>1.00</b>	<b>1.00</b>	<b>0.99</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_3$	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
(1.b)	$d_1$	<b>0.96</b>	<b>0.99</b>	<b>0.94</b>	1.00	1.00	0.98	1.00	<b>0.98</b>
	$d_2$	<b>0.98</b>	<b>1.00</b>	<b>0.96</b>	1.00	1.00	0.99	1.00	<b>0.99</b>
	$d_3$	<b>0.99</b>	<b>1.00</b>	<b>0.97</b>	1.00	1.00	0.99	1.00	<b>0.99</b>
(1.c)	$d_1$	<b>0.01</b>	<b>0.76</b>	<b>0.25</b>	0.98	0.99	0.99	1.00	<b>0.97</b>
	$d_2$	<b>0.01</b>	<b>0.90</b>	<b>0.47</b>	0.99	0.99	1.00	1.00	<b>0.98</b>
	$d_3$	<b>0.03</b>	<b>0.94</b>	<b>0.60</b>	0.99	0.99	1.00	1.00	<b>0.99</b>
(1.d)	$d_1$	<b>0.03</b>	<b>0.14</b>	<b>0.03</b>	1.00	1.00	0.82	0.56	<b>0.41</b>
	$d_2$	<b>0.06</b>	<b>0.30</b>	<b>0.06</b>	1.00	1.00	0.92	0.64	<b>0.57</b>
	$d_3$	<b>0.08</b>	<b>0.43</b>	<b>0.08</b>	1.00	1.00	0.94	0.68	<b>0.62</b>

**Table 2 :** The Proportions of  $\mathcal{P}_s$  and  $\mathcal{P}_a$ ,  $p = 2000$   $\sigma = 0.8$

Model	Size	SIRS	DC-SIS	Q-SIS	SEVIS				
		$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$				$\mathcal{P}_a$
		<i>ALL</i>	<i>ALL</i>	<i>ALL</i>	$X_1$	$X_2$	$X_{12}$	$X_{22}$	<i>ALL</i>
(1.a)	$d_1$	<b>1.00</b>	<b>1.00</b>	<b>0.97</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_2$	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_3$	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
(1.b)	$d_1$	<b>0.45</b>	<b>0.74</b>	<b>0.46</b>	1.00	1.00	0.64	1.00	<b>0.64</b>
	$d_2$	<b>0.60</b>	<b>0.85</b>	<b>0.61</b>	1.00	1.00	0.76	1.00	<b>0.76</b>
	$d_3$	<b>0.67</b>	<b>0.89</b>	<b>0.67</b>	1.00	1.00	0.83	1.00	<b>0.83</b>
(1.c)	$d_1$	<b>0.02</b>	<b>0.99</b>	<b>0.65</b>	1.00	1.00	0.88	1.00	<b>0.88</b>
	$d_2$	<b>0.04</b>	<b>1.00</b>	<b>0.80</b>	1.00	1.00	0.95	1.00	<b>0.95</b>
	$d_3$	<b>0.08</b>	<b>1.00</b>	<b>0.88</b>	1.00	1.00	0.97	1.00	<b>0.97</b>
(1.d)	$d_1$	<b>0.05</b>	<b>0.13</b>	<b>0.02</b>	1.00	0.99	0.62	0.53	<b>0.28</b>
	$d_2$	<b>0.10</b>	<b>0.27</b>	<b>0.07</b>	1.00	1.00	0.75	0.61	<b>0.42</b>
	$d_3$	<b>0.15</b>	<b>0.40</b>	<b>0.12</b>	1.00	1.00	0.82	0.68	<b>0.53</b>

**Table 3 :** The Proportions of  $\mathcal{P}_s$  and  $\mathcal{P}_a$ ,  $p = 5000$   $\sigma = 0.5$

Model	Size	SIRS	DC-SIS	Q-SIS	SEVIS				
		$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$				$\mathcal{P}_a$
		<i>ALL</i>	<i>ALL</i>	<i>ALL</i>	$X_1$	$X_2$	$X_{12}$	$X_{22}$	<i>ALL</i>
(1.a)	$d_1$	<b>1.00</b>	<b>1.00</b>	<b>0.93</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_2$	<b>1.00</b>	<b>1.00</b>	<b>0.96</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_3$	<b>1.00</b>	<b>1.00</b>	<b>0.98</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
(1.b)	$d_1$	<b>0.98</b>	<b>0.99</b>	<b>0.92</b>	1.00	1.00	0.97	1.00	<b>0.97</b>
	$d_2$	<b>0.99</b>	<b>1.00</b>	<b>0.95</b>	1.00	1.00	0.99	1.00	<b>0.99</b>
	$d_3$	<b>1.00</b>	<b>1.00</b>	<b>0.96</b>	1.00	1.00	0.99	1.00	<b>0.99</b>
(1.c)	$d_1$	<b>0.00</b>	<b>0.50</b>	<b>0.08</b>	0.98	0.98	0.96	1.00	<b>0.92</b>
	$d_2$	<b>0.01</b>	<b>0.68</b>	<b>0.19</b>	0.98	0.99	0.99	1.00	<b>0.97</b>
	$d_3$	<b>0.01</b>	<b>0.77</b>	<b>0.28</b>	0.99	0.99	0.99	1.00	<b>0.98</b>
(1.d)	$d_1$	<b>0.01</b>	<b>0.08</b>	<b>0.01</b>	0.98	0.95	0.71	0.48	<b>0.29</b>
	$d_2$	<b>0.02</b>	<b>0.18</b>	<b>0.02</b>	1.00	0.98	0.83	0.56	<b>0.43</b>
	$d_3$	<b>0.04</b>	<b>0.24</b>	<b>0.03</b>	1.00	0.98	0.88	0.60	<b>0.50</b>

**Table 4 :** The Proportions of  $\mathcal{P}_s$  and  $\mathcal{P}_a$ ,  $p = 5000$   $\sigma = 0.8$

Model	Size	SIRS	DC-SIS	Q-SIS	SEVIS				
		$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$				$\mathcal{P}_a$
		<i>ALL</i>	<i>ALL</i>	<i>ALL</i>	$X_1$	$X_2$	$X_{12}$	$X_{22}$	<i>ALL</i>
(1.a)	$d_1$	<b>1.00</b>	<b>1.00</b>	<b>0.96</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_2$	<b>1.00</b>	<b>1.00</b>	<b>0.98</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
	$d_3$	<b>1.00</b>	<b>1.00</b>	<b>0.99</b>	1.00	1.00	1.00	1.00	<b>1.00</b>
(1.b)	$d_1$	<b>0.39</b>	<b>0.66</b>	<b>0.37</b>	1.00	1.00	0.57	1.00	<b>0.57</b>
	$d_2$	<b>0.53</b>	<b>0.76</b>	<b>0.47</b>	1.00	1.00	0.69	1.00	<b>0.69</b>
	$d_3$	<b>0.59</b>	<b>0.82</b>	<b>0.54</b>	1.00	1.00	0.76	1.00	<b>0.76</b>
(1.c)	$d_1$	<b>0.01</b>	<b>0.98</b>	<b>0.50</b>	1.00	1.00	0.73	1.00	<b>0.73</b>
	$d_2$	<b>0.02</b>	<b>0.99</b>	<b>0.67</b>	1.00	1.00	0.83	1.00	<b>0.83</b>
	$d_3$	<b>0.03</b>	<b>0.99</b>	<b>0.76</b>	1.00	1.00	0.90	1.00	<b>0.90</b>
(1.d)	$d_1$	<b>0.03</b>	<b>0.08</b>	<b>0.00</b>	0.98	0.97	0.46	0.44	<b>0.14</b>
	$d_2$	<b>0.06</b>	<b>0.15</b>	<b>0.02</b>	0.99	0.99	0.62	0.55	<b>0.30</b>
	$d_3$	<b>0.08</b>	<b>0.19</b>	<b>0.04</b>	1.00	1.00	0.69	0.60	<b>0.37</b>

## Example 2

$$(2.a) \quad Y = 5X_1 + 3X_2^2 + X_3^3 + \varepsilon,$$

$$(2.b) \quad Y = 3\sin(X_1) + 7\cos(X_2) + \tan(X_3) + \varepsilon.$$

$X = (X_1, \dots, X_p)$  from a normal distribution with mean zero and covariance matrix  $\Sigma = (\sigma_{ij})_{p \times p}$ , where  $\sigma_{ii} = 1, i = 1, \dots, p$  and  $\sigma_{ij} = \sigma, i \neq j$ . Model (2.b) is a summation of trigonometric functions of predictors. Considering the property of tangent function, we generate  $X = (X_1, \dots, X_p)$  from a uniform distribution  $[-1, 1]$  and keep the covariance matrix unchanged. The random error  $\varepsilon$  follows a standard normal.

n	p	$\sigma$	d	SIRS	DC-SIS	Q-SIS	SEVIS			
				$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$			$\mathcal{P}_a$ (2.a)
				ALL	ALL	ALL	$X_1$	$X_2$	$X_3$	ALL
200	1000	0	$d_1$	<b>0.07</b>	<b>1.00</b>	<b>0.70</b>	1.00	1.00	1.00	<b>1.00</b>
200	1000	0	$d_2$	<b>0.13</b>	<b>1.00</b>	<b>0.83</b>	1.00	1.00	1.00	<b>1.00</b>
200	1000	0	$d_3$	<b>0.18</b>	<b>1.00</b>	<b>0.87</b>	1.00	1.00	1.00	<b>1.00</b>
100	1000	0	$d_1$	<b>0.04</b>	<b>0.63</b>	<b>0.15</b>	1.00	0.96	0.90	<b>0.86</b>
100	1000	0	$d_2$	<b>0.08</b>	<b>0.80</b>	<b>0.28</b>	1.00	0.98	0.93	<b>0.90</b>
100	1000	0	$d_3$	<b>0.11</b>	<b>0.85</b>	<b>0.39</b>	1.00	0.99	0.95	<b>0.93</b>
200	5000	0	$d_1$	<b>0.02</b>	<b>0.98</b>	<b>0.38</b>	1.00	1.00	0.99	<b>0.99</b>
200	5000	0	$d_2$	<b>0.03</b>	<b>0.99</b>	<b>0.54</b>	1.00	1.00	1.00	<b>1.00</b>
200	5000	0	$d_3$	<b>0.04</b>	<b>1.00</b>	<b>0.65</b>	1.00	1.00	1.00	<b>1.00</b>
100	5000	0	$d_1$	<b>0.01</b>	<b>0.32</b>	<b>0.02</b>	0.97	0.89	0.77	<b>0.65</b>
100	5000	0	$d_2$	<b>0.02</b>	<b>0.48</b>	<b>0.04</b>	0.99	0.93	0.83	<b>0.75</b>
100	5000	0	$d_3$	<b>0.02</b>	<b>0.59</b>	<b>0.08</b>	1.00	0.95	0.86	<b>0.81</b>
200	1000	0.5	$d_1$	<b>0.01</b>	<b>0.32</b>	<b>0.22</b>	1.00	0.91	0.98	<b>0.89</b>
200	1000	0.5	$d_2$	<b>0.01</b>	<b>0.47</b>	<b>0.34</b>	1.00	0.95	0.98	<b>0.93</b>
200	1000	0.5	$d_3$	<b>0.01</b>	<b>0.55</b>	<b>0.40</b>	1.00	0.96	0.99	<b>0.95</b>
100	5000	0.5	$d_1$	<b>0.00</b>	<b>0.11</b>	<b>0.03</b>	0.96	0.65	0.78	<b>0.44</b>
100	5000	0.5	$d_2$	<b>0.01</b>	<b>0.19</b>	<b>0.07</b>	0.97	0.72	0.86	<b>0.57</b>
100	5000	0.5	$d_3$	<b>0.02</b>	<b>0.26</b>	<b>0.12</b>	0.97	0.76	0.89	<b>0.64</b>
200	5000	0.5	$d_1$	<b>0.00</b>	<b>0.13</b>	<b>0.05</b>	0.99	0.79	0.94	<b>0.73</b>
200	5000	0.5	$d_2$	<b>0.00</b>	<b>0.21</b>	<b>0.09</b>	0.99	0.85	0.97	<b>0.81</b>
200	5000	0.5	$d_3$	<b>0.01</b>	<b>0.24</b>	<b>0.13</b>	1.00	0.89	0.98	<b>0.86</b>
100	5000	0.5	$d_1$	<b>0.00</b>	<b>0.01</b>	<b>0.01</b>	0.90	0.44	0.58	<b>0.19</b>
100	5000	0.5	$d_2$	<b>0.00</b>	<b>0.04</b>	<b>0.01</b>	0.92	0.52	0.65	<b>0.26</b>
100	5000	0.5	$d_3$	<b>0.00</b>	<b>0.05</b>	<b>0.03</b>	0.93	0.57	0.70	<b>0.32</b>

n	p	$\sigma$	d	SIRS	DC-SIS	Q-SIS	SEVIS			
				$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$			$\mathcal{P}_a$ (2.b)
				ALL	ALL	ALL	$X_1$	$X_2$	$X_3$	ALL
200	1000	0	$d_1$	<b>0.03</b>	<b>0.97</b>	<b>0.52</b>	1.00	1.00	1.00	<b>1.00</b>
200	1000	0	$d_2$	<b>0.07</b>	<b>0.99</b>	<b>0.73</b>	1.00	1.00	1.00	<b>1.00</b>
200	1000	0	$d_3$	<b>0.10</b>	<b>1.00</b>	<b>0.84</b>	1.00	1.00	1.00	<b>1.00</b>
100	1000	0	$d_1$	<b>0.01</b>	<b>0.27</b>	<b>0.03</b>	1.00	0.94	0.82	<b>0.77</b>
100	1000	0	$d_2$	<b>0.02</b>	<b>0.50</b>	<b>0.10</b>	1.00	0.97	0.89	<b>0.86</b>
100	1000	0	$d_3$	<b>0.04</b>	<b>0.64</b>	<b>0.20</b>	1.00	0.98	0.92	<b>0.90</b>
200	5000	0	$d_1$	<b>0.00</b>	<b>0.67</b>	<b>0.08</b>	1.00	1.00	0.95	<b>0.95</b>
200	5000	0	$d_2$	<b>0.00</b>	<b>0.82</b>	<b>0.20</b>	1.00	1.00	0.98	<b>0.98</b>
200	5000	0	$d_3$	<b>0.01</b>	<b>0.88</b>	<b>0.35</b>	1.00	1.00	0.99	<b>0.99</b>
100	5000	0	$d_1$	<b>0.00</b>	<b>0.05</b>	<b>0.00</b>	1.00	0.84	0.63	<b>0.51</b>
100	5000	0	$d_2$	<b>0.00</b>	<b>0.12</b>	<b>0.00</b>	1.00	0.91	0.72	<b>0.64</b>
100	5000	0	$d_3$	<b>0.00</b>	<b>0.19</b>	<b>0.02</b>	1.00	0.93	0.77	<b>0.71</b>
200	1000	0.5	$d_1$	<b>0.00</b>	<b>0.27</b>	<b>0.00</b>	0.98	0.93	0.98	<b>0.92</b>
200	1000	0.5	$d_2$	<b>0.02</b>	<b>0.39</b>	<b>0.00</b>	0.98	0.95	0.99	<b>0.95</b>
200	1000	0.5	$d_3$	<b>0.05</b>	<b>0.52</b>	<b>0.00</b>	0.98	0.96	0.99	<b>0.96</b>
100	1000	0.5	$d_1$	<b>0.00</b>	<b>0.06</b>	<b>0.00</b>	0.97	0.53	0.77	<b>0.41</b>
100	1000	0.5	$d_2$	<b>0.02</b>	<b>0.11</b>	<b>0.00</b>	0.98	0.64	0.85	<b>0.53</b>
100	1000	0.5	$d_3$	<b>0.03</b>	<b>0.16</b>	<b>0.00</b>	0.98	0.72	0.90	<b>0.64</b>
200	5000	0.5	$d_1$	<b>0.00</b>	<b>0.09</b>	<b>0.00</b>	0.98	0.81	0.94	<b>0.77</b>
200	5000	0.5	$d_2$	<b>0.00</b>	<b>0.14</b>	<b>0.00</b>	0.98	0.87	0.96	<b>0.85</b>
200	5000	0.5	$d_3$	<b>0.00</b>	<b>0.17</b>	<b>0.00</b>	0.98	0.90	0.97	<b>0.88</b>
100	5000	0.5	$d_1$	<b>0.00</b>	<b>0.01</b>	<b>0.00</b>	0.97	0.36	0.59	<b>0.20</b>
100	5000	0.5	$d_2$	<b>0.00</b>	<b>0.02</b>	<b>0.00</b>	0.98	0.44	0.68	<b>0.29</b>
100	5000	0.5	$d_3$	<b>0.00</b>	<b>0.03</b>	<b>0.00</b>	0.98	0.50	0.72	<b>0.35</b>



### Example 3

The underlying model is adapted from Ravikumar *et al.* (2008). We generate the data from the additive model

$$(3) \quad Y = f_1(X_1) + f_2(X_2) + f_3(X_3) + f_4(X_4) + \varepsilon,$$

where

$$f_1(x) = -2\sin(2x), \quad f_2(x) = x^2 - 1/3, \quad f_3(x) = x - 1/2, \quad f_4(x) = e^{-x} + e^{-1} - 1.$$

$X_1, \dots, X_p$  are drawn from an independent and identically distributed uniform distribution on  $[-2, 2]$ . The random error  $\varepsilon$  follows the standard normal distribution. Let  $n = 150$  and  $p = 2000$  which is ten times of that in *et al.* (2008). The censored observations  $Y_i^* = \min(Y_i, C_i)$ ,  $1 \leq i \leq n$ , where  $C_i$  follows a uniform distribution on  $[0, 8.5]$  to control the censoring rate to be approximately 25%. The replications are 500.

**Table 5 :** The Proportions of  $\mathcal{P}_s$  and  $\mathcal{P}_a$  in Example 3

d	SIRS	DC-SIS	Q-SIS	SEVIS				
	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_a$	$\mathcal{P}_s$				$\mathcal{P}_a$
	<i>ALL</i>	<i>ALL</i>	<i>ALL</i>	$X_1$	$X_2$	$X_3$	$X_4$	<i>ALL</i>
$d_1$	<b>0.00</b>	<b>0.16</b>	<b>0.01</b>	0.93	0.90	0.97	1.00	<b>0.81</b>
$d_2$	<b>0.00</b>	<b>0.33</b>	<b>0.02</b>	0.97	0.94	0.98	1.00	<b>0.89</b>
$d_3$	<b>0.00</b>	<b>0.46</b>	<b>0.05</b>	0.98	0.96	0.99	1.00	<b>0.92</b>

# Outline

- 1 What've been done in independence screening
- 2 Sure Explained Variability: Definitions and Properties
- 3 Sure explained variability and independence screening (SEVIS)
- 4 Illustrative exmaples
- 5 Integrated genomic explained variability analyses of ovarian carcinoma**
- 6 A nonparametric estimator for SEV

## Data

- 593 patients and their 12042 gene expression levels
- Divide all patients into 17 groups by tissue source site (TSS)
- 8 groups: TCGA-04, TCGA-09, TCGA-13, TCGA-23, TCGA-24, TCGA-25, TCGA-29, TCGA-61
- At least 30 patients ( $n \geq 30$ ), total 487 samples
- Response variable 1: TP53, a well known tumor suppressor gene and mutates in 303 among 316 ovarian carcinoma patients
- Response variable 2: survival time of patients

## Evaluate the performance of those top ranked genes

- consider the following nonparametric additive model,

$$Y = f_1(X_{first}) + f_2(X_{second}) + f_3(X_{third}) + \varepsilon, \quad (5.1)$$

- compare the adjusted  $R^2$  and deviance explained by each method
- use  $SEV(Y|\hat{Y})$  to measure nonlinearly explained variability, where  $\hat{Y} = \hat{f}_1(X_{first}) + \hat{f}_2(X_{second}) + \hat{f}_3(X_{third})$

Table 6 : The results of data fitting for  $Y = \text{TP53}$

Group	n	SIRS			DC-SIS		
		adj $R^2$	deviance explained	SEV( $Y Y$ )	adj $R^2$	deviance explained	SEV( $Y Y$ )
TCGA-04	43	0.405	0.467	0.410	0.578	0.641	0.562
TCGA-09	30	0.451	0.529	0.435	0.629	0.697	0.577
TCGA-13	113	0.294	0.342	0.311	0.319	0.373	0.335
TCGA-23	38	0.450	0.495	0.419	0.534	0.635	0.556
TCGA-24	100	0.296	0.361	0.326	<b>0.311</b>	<b>0.373</b>	<b>0.337</b>
TCGA-25	45	0.310	0.359	0.315	0.704	0.773	0.652
TCGA-29	52	0.556	0.583	0.525	0.539	0.566	0.484
TCGA-61	66	<b>0.494</b>	<b>0.558</b>	<b>0.489</b>	0.440	0.505	0.455

**Table 7 :** The results of data fitting for  $Y = \text{TP53}$

Group	n	Q-SIS			SEVIS		
		adj $R^2$	deviance explained	SEV( $Y Y$ )	adj $R^2$	deviance explained	SEV( $Y Y$ )
TCGA-04	43	0.303	0.447	0.404	<b>0.861</b>	<b>0.912</b>	<b>0.788</b>
TCGA-09	30	0.200	0.400	0.341	<b>0.677</b>	<b>0.774</b>	<b>0.654</b>
TCGA-13	113	0.188	0.249	0.243	<b>0.348</b>	<b>0.421</b>	<b>0.382</b>
TCGA-23	38	0.538	0.681	0.591	<b>0.625</b>	<b>0.717</b>	<b>0.684</b>
TCGA-24	100	0.181	0.250	0.237	0.230	0.294	0.278
TCGA-25	45	0.562	0.716	0.613	<b>0.906</b>	<b>0.956</b>	<b>0.815</b>
TCGA-29	52	0.530	0.585	0.505	<b>0.620</b>	<b>0.643</b>	<b>0.545</b>
TCGA-61	66	0.151	0.246	0.272	0.267	0.365	0.380

**Table 8 :** The results of data fitting for  $Y$  being days to last followup

Group	n	SIRS			DC-SIS		
		adj $R^2$	deviance explained	SEV( $Y Y$ )	adj $R^2$	deviance explained	SEV( $Y Y$ )
TCGA-04	38	0.822	0.894	0.777	<b>0.884</b>	<b>0.940</b>	<b>0.799</b>
TCGA-09	30	0.527	0.577	0.507	0.720	0.778	0.646
TCGA-13	109	0.274	0.294	0.265	0.287	0.361	0.334
TCGA-23	38	0.660	0.725	0.639	0.662	0.720	0.657
TCGA-24	98	0.321	0.364	0.371	0.258	0.293	0.273
TCGA-25	45	0.456	0.493	0.456	<b>0.573</b>	<b>0.660</b>	<b>0.589</b>
TCGA-29	49	0.401	0.439	0.389	<b>0.503</b>	<b>0.554</b>	0.483
TCGA-61	64	0.224	0.261	0.241	0.302	0.336	0.311



**Table 9 :** The results of data fitting for  $Y$  being days to last followup

Group	n	Q-SIS			SEVIS		
		adj $R^2$	deviance explained	SEV( $Y Y$ )	adj $R^2$	deviance explained	SEV( $Y Y$ )
TCGA-04	38	0.655	0.792	0.666	0.710	0.794	0.740
TCGA-09	30	0.533	0.641	0.597	<b>0.725</b>	<b>0.800</b>	<b>0.728</b>
TCGA-13	109	0.257	0.301	0.278	<b>0.355</b>	<b>0.403</b>	<b>0.431</b>
TCGA-23	38	0.476	0.531	0.523	<b>0.691</b>	<b>0.752</b>	<b>0.690</b>
TCGA-24	98	0.289	0.361	0.382	<b>0.330</b>	<b>0.400</b>	<b>0.423</b>
TCGA-25	45	0.317	0.364	0.329	0.528	0.622	0.547
TCGA-29	49	0.320	0.414	0.364	0.398	0.495	<b>0.514</b>
TCGA-61	64	0.253	0.345	0.342	<b>0.332</b>	<b>0.416</b>	<b>0.429</b>

## the mutual pairwise explained variabilities

- the explained variabilities of more than 72.5% of genes selected using SIRS, DC-SIS and Q-SIS by the genes selected using SEVIS are larger than the explained variabilities in the opposite directions with the response variable being TP53.
- The proportion increases to 81.4% with the response variable being changed to days to last followup.

# Outline

- 1 What've been done in independence screening
- 2 Sure Explained Variability: Definitions and Properties
- 3 Sure explained variability and independence screening (SEVIS)
- 4 Illustrative examples
- 5 Integrated genomic explained variability analyses of ovarian carcinoma
- 6 A nonparametric estimator for SEV**

- Suppose  $f(x, y)$  is the joint density function of random variables  $(X, Y)$ ,
- and  $f^X(x)$  is the marginal density function of predictor  $X$ .

$$\begin{aligned} \text{SEV}(Y|X) &= \frac{E[E(Y|X)]^2 - (E(Y))^2}{\text{var}(Y)} \\ &= \frac{\int (r^{Y|X}(x))^2 f^X(x) dx - (E(Y))^2}{\text{var}(Y)}, \end{aligned} \quad (6.1)$$

where

$$r^{Y|X}(x) := E(Y|X=x) := \frac{\phi^{Y|X}(x)}{f^X(x)} := \frac{\int y f(x, y) dy}{f^X(x)}. \quad (6.2)$$

## Nadaraya-Watson estimators

$$\begin{aligned}\phi_n^{Y|X}(x) &= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) Y_i, \\ f_n^X(x) &= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right),\end{aligned}\tag{6.3}$$

$$r_n^{Y|X}(x) = \begin{cases} \frac{\sum_{i=1}^n K(\frac{x-X_i}{h}) Y_i}{\sum_{j=1}^n K(\frac{x-X_j}{h})} & \text{if } \sum_{j=1}^n K(\frac{x-X_j}{h}) \neq 0, \\ 0 & \text{if } \sum_{j=1}^n K(\frac{x-X_j}{h}) = 0. \end{cases}\tag{6.4}$$

The estimator

$$\hat{\omega} = \widehat{\text{SEV}}(Y|X) = S_Y^{-2} \left( \int (r_n^{Y|X}(x))^2 f_n^X(x) dx - (\bar{Y})^2 \right),\tag{6.5}$$

## Nadaraya-Watson estimators

$$\begin{aligned}\phi_n^{Y|X}(x) &= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) Y_i, \\ f_n^X(x) &= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right),\end{aligned}\tag{6.3}$$

$$r_n^{Y|X}(x) = \begin{cases} \frac{\sum_{i=1}^n K(\frac{x-X_i}{h}) Y_i}{\sum_{j=1}^n K(\frac{x-X_j}{h})} & \text{if } \sum_{j=1}^n K(\frac{x-X_j}{h}) \neq 0, \\ 0 & \text{if } \sum_{j=1}^n K(\frac{x-X_j}{h}) = 0. \end{cases}\tag{6.4}$$

## The estimator

$$\hat{\omega} = \widehat{\text{SEV}}(Y|X) = S_Y^{-2} \left( \int (r_n^{Y|X}(x))^2 f_n^X(x) dx - (\bar{Y})^2 \right),\tag{6.5}$$

## Asymptotics

Under Assumptions (A1) and (A2), we have

$$\sqrt{n}(\mathbf{A}^T \Sigma \mathbf{A})^{-\frac{1}{2}}(\widehat{\text{SEV}}(Y|X) - \text{SEV}(Y|X)) \Rightarrow N(0, 1), \quad (6.6)$$

where

$$\Sigma = \text{Cov}\left(\sigma_Y^{-2} \int (2Y_i - \frac{\phi^{Y|X}(x)}{f^X(x)}) \frac{\phi^{Y|X}(x)}{f^X(x)} \frac{1}{h} K(\frac{x - X_i}{h}) dx, \frac{Y_i}{\sigma_Y}, Y_i^2\right),$$

and

$$\mathbf{A} = (1, -\frac{2\mu_Y}{\sigma_Y} + 2\mu_Y \sigma_Y^{-3} (\int \frac{(\phi^{Y|X}(x))^2}{f^X(x)} dx - \mu_Y^2), -\sigma_Y^{-4} (\int \frac{(\phi^{Y|X}(x))^2}{f^X(x)} dx$$

$\mu_Y$  and  $\sigma_Y$  are the population mean and standard deviation of  $Y$ , respectively.