

Outline: Ultra-High Dimensional Supervised Problems (F. Chiaromonte)

Statistical Methods for Large, Complex Data

- Main reference: Fan J., Lv J. (2008) Sure Independence Screening for ultrahigh dimensional feature spaces. JRSS-B, 70(5) 849-911.
- Additional references: through the slides.

Stochastic mechanism:

Y, ε r.vbls in R^1 X r.vct in R^p

$$E(X) = 0 \quad Cov(X) = \Sigma_X$$

$$\varepsilon \text{ indep } X \quad \varepsilon \sim N(0, \sigma^2)$$

Independent sampling:

$$Y_{nx1} = X_{nxp} \beta_{px1} + \varepsilon_{nx1}$$

Estimation (fitting)

$$\hat{\beta}_{LS} = (X'X)^{-1}X'Y$$

- Take the marginal (linear) associations between Y and each of the X's; $X'Y$
- “Re-map” it based on the (linear) associations among the X's; S_x .

How would you rank the X's and chose the d << p most likely to have an effect on Y?

On the basis of the effects as estimated by LS; rank the coordinates of the LS vector and pick the largest d.

Does this work?

Yes, if the re-mapping is effective; the ranking comprises two ingredients, marginal associations with Y and re-mapping based on associations among the X's. No problem at all if “no re-mapping” i.e. S_x proportional to I_p .

What can go wrong?

$\lambda_{\max}(S_x)$ is large and $\lambda_{\min}(S_x)$ is small...

- Linear “concentration” of the data cloud in feature space, collinearity in the sample. Because Σ_x contains collinearity, and/or because n is not large enough wrt p .
- An ineffective re-mapping inflates the sampling variability in the LS estimates of the feature effects, possibly makes them poorly determined or non-unique on any given sample.

What happens as p grows wrt n ; $p \sim n$, $p > n$, $p \gg n$?

- Σ_x is likely to comprise more collinearity; we are considering more and more features which may carry linear associations, and
- S_x becomes progressively more collinear than its population counterpart, just because we have insufficient replication

What do we do?

$$\hat{\beta}_{MP} = (X'X)^+ X'Y \quad \text{More-Penrose generalized inverse}$$

$$\hat{\beta}_{Ridge}(\lambda) = (X'X + \lambda I)^{-1} X'Y \quad \text{Regularize with a size constraint, L2}$$

$$\hat{\beta}_{LASSO}(\lambda) = \operatorname{argmin} \{ \|Y - X\beta\|^2 + \lambda L1(\beta) \} \quad \text{Regularize and sparsify with a size constraint, L1}$$

... and more sophisticated approaches; SCAD, Adaptive LASSO, Danzig Selector. Also old fashioned Partial Least Squares.

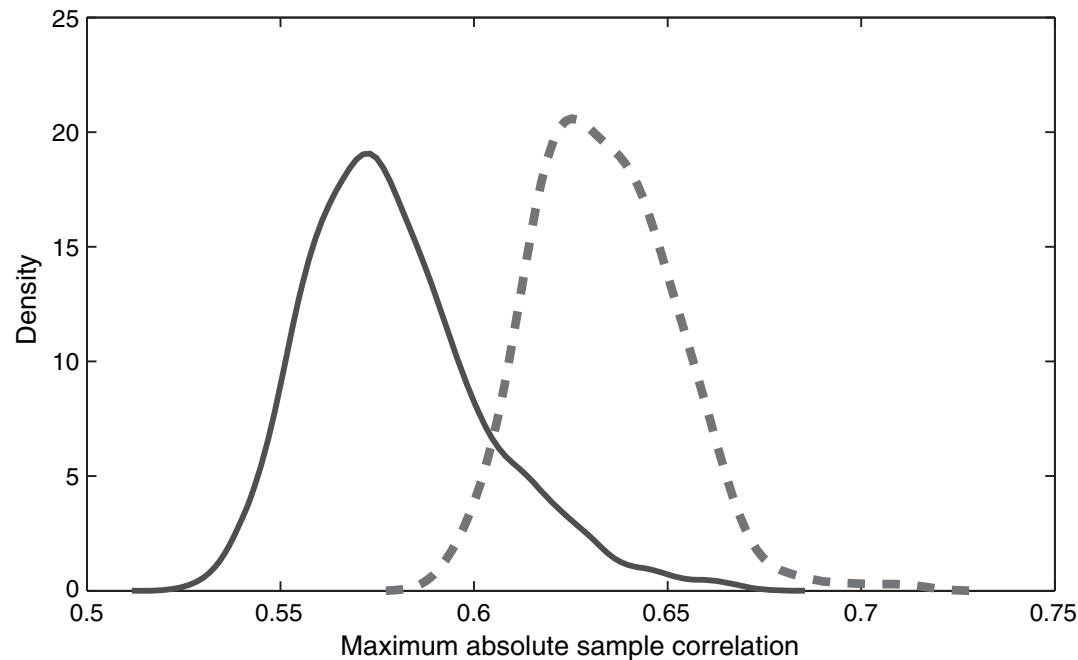


Fig. 1. Distributions of the maximum absolute sample correlation coefficient when $n = 60$ and $p = 1000$ (—) and $n = 60$ and $p = 5000$ (- - -), based on 500 simulations

X data simulated with Σ_x proportional to I_p .

No correlations at the population level, yet if $p \gg n$ the correlations in S_x are large!

What if we ignored the re-mapping all together?

Base the ranking and choice of d features on the component-wise regression, i.e., on the marginal correlations

$$\hat{\beta}_{LS} = (X'X)^{-1}X'Y$$

$$\hat{\beta}_{LS} = \omega \text{ if } X'X \propto I_p$$

$$\hat{\beta}_{Ridge}(\lambda) = (X'X + \lambda I)^{-1}X'Y$$

$$\begin{aligned}\hat{\beta}_{Ridge}(\lambda) &\xrightarrow{\lambda \rightarrow \infty} 0 \\ \lambda \hat{\beta}_{Ridge}(\lambda) &\xrightarrow{\lambda \rightarrow \infty} \omega\end{aligned}$$

$$\omega = X'Y \propto r_{yx}$$

CORRELATION LEARNING

as the penalty increases the Ridge vector shrinks to 0, but if we multiply it by λ (which does not change the ranking of the components) ... it converges to ω !!

More generally, approaches to learn about importance of features **marginally**; these can be applied with small computational burden and without the curse of dimensionality even when $p \gg n$

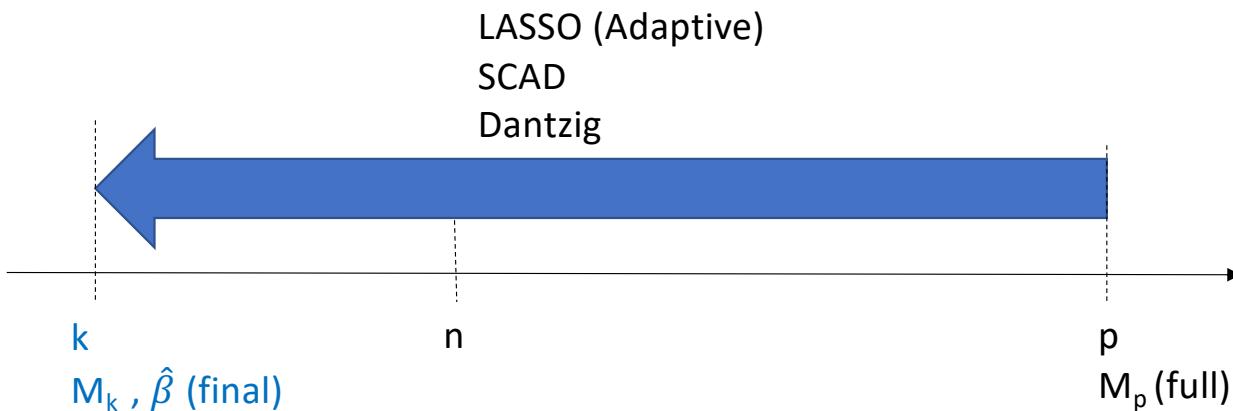
IMPORTANT REMARK: A NEW ASYMPTOTICS

In addition to problems with $p \gg n$, starting with the new century we have problems where $p = p(n)$. As n grows, p grows with it, possibly exponentially $\log(p) = O(n^a)$ though perhaps with small a .

The more we sample, observing more units, the more the dimension of our feature space increases. Think of variants (SNPs) as we sequence genomes of more individuals, or recorded product choices/scores as we branch out over a network of individuals on social media. The "universe" of the observed SNPs (of products) increases as we include more people.

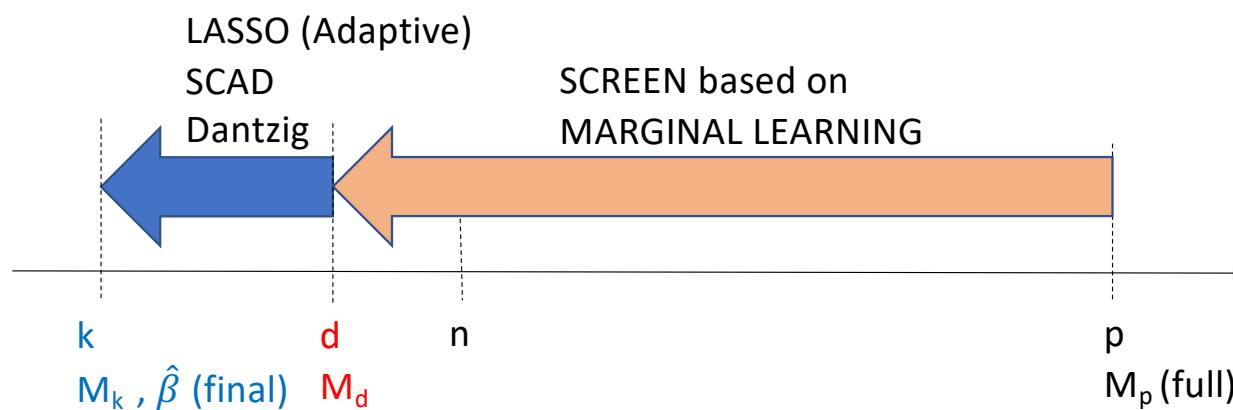
In Statistics this required the formulation of a **new type of asymptotics**; study the performance of statistical procedures as n grows, and $p(n)$ grows with it – possibly not too fast.

The proposal of Fan and Lv (2008) : If we believe the underlying true model is sparse, instead of



Screening disregards associations among the X's, but when $p \gg n$ our ability to account for them and re-map is so poor that it may hurt more than it helps.

use a reasonable $d < n$, e.g., $d = n-1$, or $d = n/\log(n)$ and proceed in two stages



To get to the true sparse model M_* with s features, and to estimate its true β , the two-stage strategy (which uses only marginal information when the dimension is still ultra-high) is more effective!!

Not a new idea, practitioners have used it forever, but now formalized. And its performance established by simulations and “new asymptotics” theoretical results.

SURE INDEPENDENCE SCREENING (SIS)

- SURE SCREENING PROPERTY: with a reasonable $d = n-1$ or $n/\log(n)$ and some assumptions on the nature of the stochastic mechanism generating the data and the speed at which $p(n) \gg n$ grows with n , when n is large $\Pr(M_* \text{ is contained in } M_d)$ is “overwhelming” (growing to 1).
- INDEPENDENCE: refers to the fact that we operate on **marginal information**, one X at a time.

The SCREENING algorithm is trivial and computationally inexpensive:

- Compute the marginal correlations $\omega = X'Y \propto r_{yx}$
- Rank the entries in ω and pick the first d to form M_d .

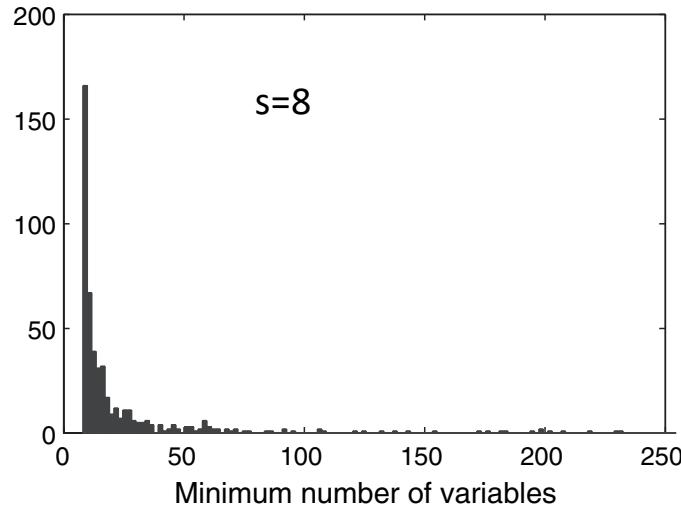
Next, apply feature selection algorithm XXX to the regression of Y on the features contained in M_d to obtain $M_k, \hat{\beta}$ (final).

First, use **simulations** to assess whether, on finite (but large) samples

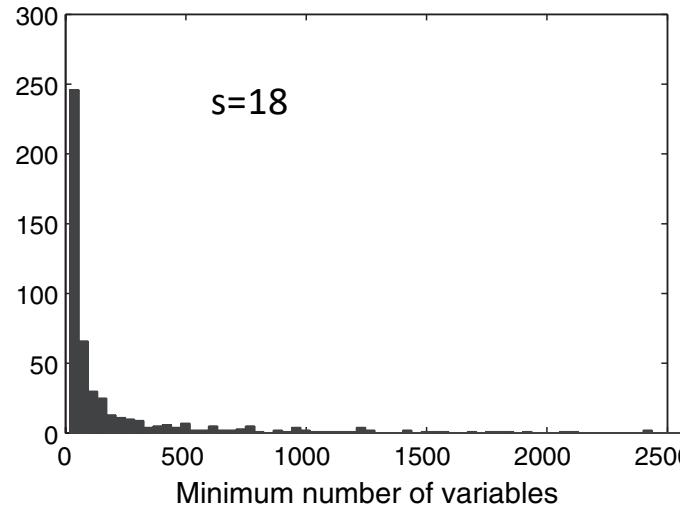
- the sure screening property holds for SIS, and
- M_k is close to M_* (k close to s) and $\hat{\beta}$ close to the true β for SIS-XXX

In particular

- do we do better than with one pass of XXX? and
- is performance affected by the presence of collinearity in Σ_x (population level)?



(a)



(b)

Fig. 5. Distribution of the minimum number of selected variables that is required to include the true model by using SIS when (a) $n = 200$ and $p = 1000$ and (b) $n = 800$ and $p = 20\,000$ in simulation I

OK!

Table 1. Results of simulation I: medians of the selected model sizes and estimation errors (in parentheses)

p	one pass		Results for the following methods:			
	Dantzig selector	Lasso	SIS-SCAD	SIS-DS	SIS-DS-SCAD	SIS-DS-AdaLasso
1000 s=8	10^3 (1.381)	62.5 (0.895)	15 (0.374)	37 (0.795)	27 (0.614)	34 (1.269)
20000 s=18	—	—	37 (0.288)	119 (0.732)	60.5 (0.372)	99 (1.014)

One pass does very poorly at imposing sparsity and at estimating effects. And for $p=20,000$ cannot be run in reasonable time.

SIMULATIONS WITH UNCORRELATED X'S

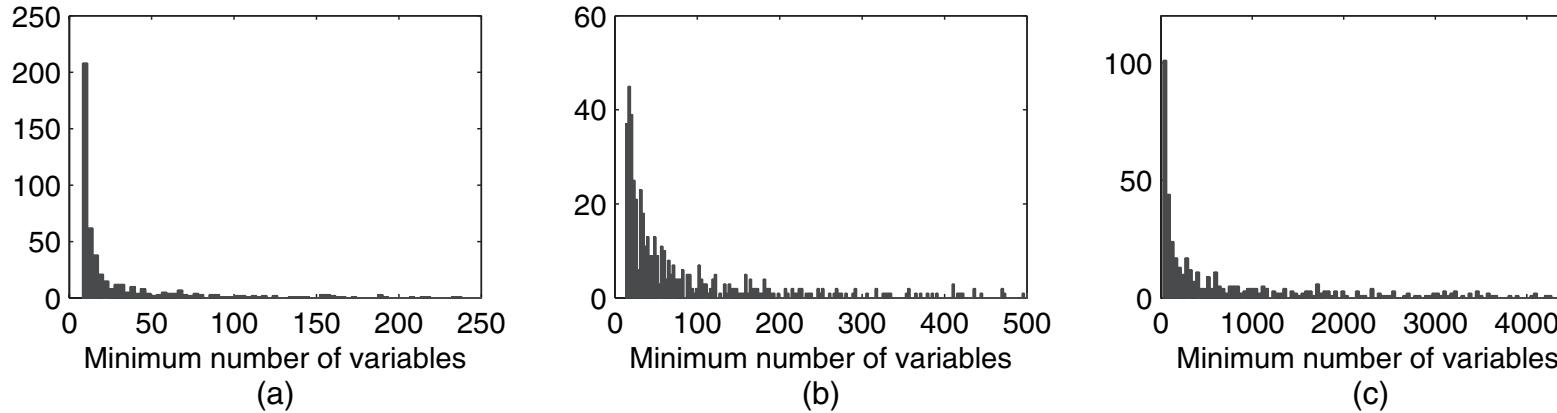


Fig. 6. Distribution of the minimum number of selected variables that is required to include the true model by using SIS when (a) $n = 200$, $p = 1000$ and $s = 5$, (b) $n = 200$, $p = 1000$ and $s = 8$ and (c) $n = 800$, $p = 20000$ and $s = 8$ in simulation II $s=14$?

Table 2. Results of simulation II: medians of the selected model sizes and estimation errors (in parentheses)

k $\|\hat{\beta} - \beta\|^2$

p	one pass		Results for the following methods:			
	Dantzig selector	Lasso	SIS-SCAD	SIS-DS	SIS-DS-SCAD	SIS-DS-AdaLasso
1000 ($s=5$)	10^3 (1.256)	91 (1.257)	21 (0.331)	56 (0.727)	27 (0.476)	52 (1.204)
	10^3 (1.465)	74 (1.257)	18 (0.458)	56 (1.014)	31.5 (0.787)	51 (1.824)
20000 $s=14$?	—	—	36 (0.367)	119 (0.986)	54 (0.743)	86 (1.762)

OK! also with correlated X's

One pass does very poorly at imposing sparsity and at estimating effects. And for $p=20,000$ cannot be run in reasonable time.

Theoretical Assessments

SURE SCREENING PROPERTY OF SIS:

- Let $M_d = \text{SIS}(M_p)$ with $d = n/\log(n)$
- *Impose assumptions on the nature of the stochastic mechanism generating the data, and on the speed at which $p(n) >> n$ grows with n*
- Prove $\Pr(M_* \subseteq M_d) \geq 1 - (\text{fast vanishing as } n \rightarrow \infty)$

Also, again under appropriate assumption, prove

- CONSISTENCY FOR SIS-DANTZIG $\hat{\beta} \xrightarrow[n \rightarrow \infty]{} \beta$
- ORACLE PROPERTY FOR SIS-SCAD
 - (i) $\hat{\beta}_j = 0$ for any $j \notin M_*$
 - (ii) $\hat{\beta}$ does as well for β as if we had run standard LS **knowing** the true M_*

Assumptions

Stochastic mechanism:

$$Y, \varepsilon \text{ r.vbls in } R^1 \quad X, Z = \Sigma_X^{-1/2} X \text{ r.vct in } R^p$$

$$E(Z) = 0 \quad Cov(Z) \propto I_p$$

$$\varepsilon \text{ indep } Z \quad \varepsilon \sim N(0, \sigma^2)$$

(I) $\text{Var}(Y)$ does not grow, and the effects for the relevant features are strong enough, as n grows; for some $\kappa \geq 0$

(II) The relevant features have a strong enough “trace” in their marginal linear associations with Y

$p > n$ grows exponentially, but slow enough relative to the strength of the signals; for some $0 < \alpha < (1-2\kappa)$

(III) The feature collinearity (population level) remains weak enough as $p(n)$ grows ; for some $\tau \geq 0$

Z is spherically symmetric and, for any submatrix with more than n columns \tilde{Z}_{nxh} , $n < h \leq p$, $\tilde{Z}_{nxh}\tilde{Z}_{nxh}'$ has all n eigenvalues of the same order – not too dissimilar (adds to the required symmetry for the features). Nothing much beyond Σ_X . This certainly holds if X is Gaussian (Z standard Gaussian).

Independent sampling:

$$Y_{nx1} = X_{nxp}\beta_{px1} + \varepsilon_{nx1} = Z_{nxp}\Sigma_X^{1/2}\beta_{px1} + \varepsilon_{nx1}$$

$$\min_{j \in M_*} |\beta_j| \geq \frac{\text{pos const}}{n^\kappa}$$

$$\min_{j \in M_*} \left| cov \left(\frac{Y}{\beta_J}, X_j \right) \right| \geq \text{pos const}$$

$$\log(p) = O(n^\alpha)$$

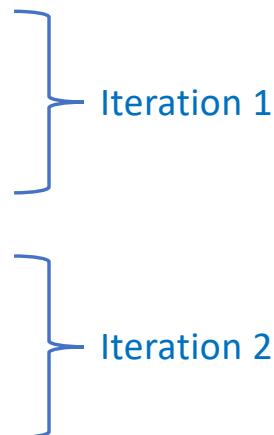
$$\lambda_{\max}(\Sigma_X) \leq \text{pos const} \times n^\tau$$

SIS-XXX can perform poorly though (even for fairly large n) if:

- The effects of some relevant features are weaker than the spurious marginal traces carried by non-relevant features correlated with strong relevant ones; failure in (I), (III).
- Some relevant features have weak marginal linear traces – they may act only jointly with other features, or act with little linear trend; failure in (II).

Performance can be rescued through another very old and effective idea: iterate the procedure, working with residuals as to magnify weaker signals: **ITERATIVE SIS**, ISIS-XXX. The algorithm:

- SIS for Y on $M_p \Rightarrow M_{dIT}^{(1)}$
- XXX for Y on $M_{dIT}^{(1)} \Rightarrow (M_{k(1)}, \hat{\beta}_{(1)})$
- $A^{(1)} = M_{k(1)}$
- $e^{(1)} = Y - X_{A(1)} \hat{\beta}_{(1)}$
- SIS for $e^{(1)}$ on $M_p \setminus A^{(1)} \Rightarrow M_{dIT}^{(2)}$
- XXX for $e^{(1)}$ on $M_{dIT}^{(2)} \Rightarrow (M_{k(2)}, \hat{\beta}_{(2)})$
- $A^{(2)} = A^{(1)} \cup M_{k(1)}$
- $e^{(2)} = Y - X_{A(2)} \hat{\beta}_{(2)}$



Iteration 1
Iteration 2

... etc. (usually just a few times). Stop iterating when $\#(A^{(h)})$ reaches $n-1$ or $n/\log(n)$

XXX for Y on $A^{(h)} \Rightarrow (M_k, \hat{\beta})$ the final result.

Table 1. Results of simulation I: medians of the selected model sizes and estimation errors (in parentheses)

<i>p</i>	Results for the following methods:					
	Dantzig selector	Lasso	SIS-SCAD	SIS-DS	SIS-DS-SCAD	SIS-DS-AdaLasso
1000 s=8	10^3 (1.381)	62.5 (0.895)	15 (0.374)	37 (0.795)	27 (0.614)	34 (1.269)
20000	—	—	37 (0.288)	119 (0.732)	60.5 (0.372)	99 (1.014)
s=18	—	—				

Table 2. Results of simulation II: medians of the selected model sizes and estimation errors (in parentheses)

<i>p</i>	Results for the following methods:					
	Dantzig selector	Lasso	SIS-SCAD	SIS-DS	SIS-DS-SCAD	SIS-DS-AdaLasso
1000 (<i>s</i> =5)	10^3 (1.256)	91 (1.257)	21 (0.331)	56 (0.727)	27 (0.476)	52 (1.204)
(<i>s</i> =8)	10^3 (1.465)	74 (1.257)	18 (0.458)	56 (1.014)	31.5 (0.787)	51 (1.824)
20000	—	—	36 (0.367)	119 (0.986)	54 (0.743)	86 (1.762)
s=14?	—	—				

Table 7. Simulations I and II in Section 3.3 revisited: medians of the model sizes selected and the estimation errors (in parentheses) for the SIS-SCAD method

<i>p</i>	Results for simulation I	Results for simulation II
1000 s=8	13 (0.329)	(<i>s</i> =5) 11 (0.223) (<i>s</i> =8) 13.5 (0.366)
20000 s=14?	31 (0.246)	s=14? 27 (0.315)

Iterations work!

In the article more simulation results to demonstrate effectiveness against more specific failure scenarios

After the seminal paper by Fan and Lv (2008), there has been a deluge of literature on screening algorithms for ultra-high dimensional supervised problems.

- **Extension to GLMs:** Fan J., Samworth R., Wu Y. (2009) Ultrahigh dimensional feature selection: beyond the linear model. *Journal of Machine Learning Research*, 10 2013-2038.
- **Extension to model-free settings:** Zhu L., Li L., Li R., Zhu L. (2011). Model-free feature screening for ultra-high dimensional data. *JASA* 106(496) 1464-1475.

... and many others. Still a very active area of research.

What changes is the implementation of marginal learning; ω is defined differently (based on a GLM framework, or non-parametrically). Also, theoretical results can be different and require alternative proof strategies.

Remark: a screening algorithm must be *conservative, the focus is on controlling false negatives*. We want to make sure we do not leave out any relevant feature; more features can be eliminated as needed by feature selection after screening. If we consider SIS-XXX as a whole, it makes sense to also assess performance in terms of false positives (we do not want to finish with a k much larger than s ; the size of the true sparse model) and in terms of estimation quality (how close $\hat{\beta}$ comes to the true β).

Some pointers on the use of covariate information
matrices/numbers for supervised dimension reduction
and feature screening

SUFFICIENT DIMENSION REDUCTION BASED ON INFORMATION MATRICES



Theory and Methods

Covariate Information Matrix for Sufficient Dimension Reduction

Weixin Yao Debmalya Nandy, Bruce G. Lindsay & Francesca Chiaromonte

Pages 1752-1764 | Received 03 Mar 2017, Accepted 10 Aug 2018, Accepted author version posted online: 06 Sep 2018, Published online: 27 Feb 2019

Download citation <https://doi.org/10.1080/01621459.2018.1515080>

Check for updates



Covariate Information Matrix (CIM): combine “local” non-parametric assessment of densities and “global” Eigen-decomposition, in an **Information Matrix framework...**

... prior applications to:

- Projection Pursuit
- Spherical Symmetry, Multivariate Structures
- Independent Components Analysis
- Graphical Models

Hui G., Lindsay B. (2010). Projection pursuit via white noise matrices . *Sankhya B*, 72(2): 123–153 .

Lindsay B., Yao W. (2012). Fisher information matrix: A tool for dimension reduction, projection pursuit, independent component analysis, and more. *Canadian Journal of Statistics*, 40(4): 712–730.

What is the theory behind this?

$$U_{\mathbf{x}}(y) = \nabla_{\mathbf{x}} \log f(y \mid \mathbf{x})$$

score vector

$$\mathbb{F}_{\mathbf{x}} = \int U_{\mathbf{x}}(y) U_{\mathbf{x}}(y)^T f(y \mid \mathbf{x}) dy$$

Fisher Information Matrix for “parameter” x

How $Y \mid X=x$ changes with x

$$\mathbb{C}_{\mathbf{X}} = \int \mathbb{F}_{\mathbf{x}} f(\mathbf{x}) d\mathbf{x}$$

average on X ... **Covariate Information Matrix**

Key result: under very general conditions one has $\text{Span}[\mathbb{C}_X] \equiv \Sigma_X S_{Y|X}$ (the Central Subspace of SDR)

$$U_f(\mathbf{x}) = \nabla_{\mathbf{x}} \log f(\mathbf{x})$$

$$\mathbb{J}_{\mathbf{X}} = \int U_f(\mathbf{x}) U_f(\mathbf{x})^T f(\mathbf{x}) d\mathbf{x}$$

Density Information Matrix for X
characterization of the X distribution

$$U_{f^{(y)}}(\mathbf{x}) = \nabla_{\mathbf{x}} \log f^{(y)}(\mathbf{x})$$

$$\mathbb{J}_{\mathbf{X}|Y=y} = \int U_{f^{(y)}}(\mathbf{x}) U_{f^{(y)}}(\mathbf{x})^T f^{(y)}(\mathbf{x}) d\mathbf{x}, \quad \text{Density Information Matrix for } X \mid Y=y$$

$$\mathbb{J}_{\mathbf{X}|Y} = \int \mathbb{J}_{\mathbf{X}|Y=y} f(y) dy.$$

average on Y ... **Inverse Regression**

Key result: $\mathbb{C}_X = \mathbb{J}_{\mathbf{X}|Y} - \mathbb{J}_{\mathbf{X}}$

rewrite CIM in terms of X -related DIMs
Inverse Regression, “adjusted” for X

What is the implementation?

(2) “global” *Eigen-decomposition* $\hat{\mathbb{C}}_X = \sum_{j=1}^p \lambda_j v_j v_j^T \rightarrow \hat{\mathcal{S}}_{Y|X} = \hat{\Sigma}_X^{-1} \text{Span}[v_1 \dots v_d]$

Key result: without (L) and (C), one has $\text{Span}[\mathbb{C}_X] \equiv \Sigma_X \mathcal{S}_{Y|X}$

(1) Inexpensive, effective estimation of $f(x)$ and $f(x|y_\ell)$, $\ell = 1 \dots L$ (slices) with the *f2-method for non-parametric density estimation* (use of squared surrogates preserving peaks and valleys) $\hat{J}_X, \hat{J}_{X|Y} \rightarrow \hat{\mathbb{C}}_X$

Key result: $\mathbb{C}_X = \mathbb{J}_{X|Y} - \mathbb{J}_X$

What is the implementation?

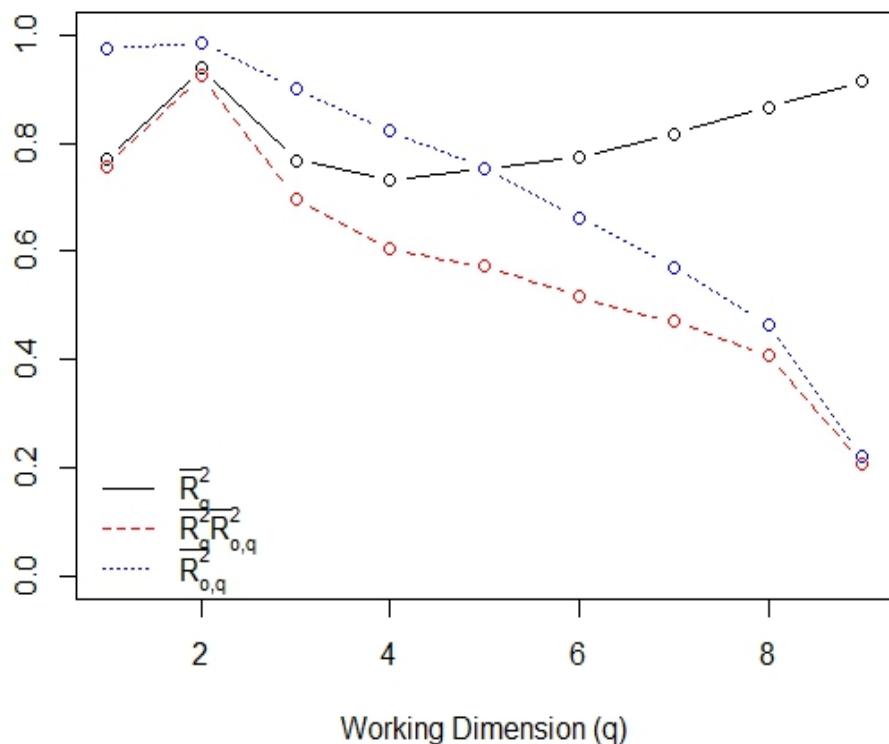
(2) “global” *Eigen-decomposition* $\hat{\mathbb{C}}_X = \sum_{j=1}^p \lambda_j v_j v_j^T \rightarrow \hat{\mathcal{S}}_{Y|X} = \hat{\Sigma}_X^{-1} \text{Span}[v_1 \dots v_d]$

dimension estimation
new bootstrap-based
diagnostic plot (see later)

(1) Inexpensive, effective estimation of $f(x)$ and $f(x|y_\ell)$, $\ell = 1 \dots L$ (slices)
with the *f2-method for non-parametric density estimation* (use of squared
surrogates preserving peaks and valleys) $\hat{J}_X, \hat{J}_{X|Y} \rightarrow \hat{\mathbb{C}}_X$

Diagnostic plot for dimension estimation

(applicable with any SDR method)



Heteroscedastic Y scenario

$\sigma = 0.2$, Independent \mathbf{X} , $n = 400$

CIM with 5 slices

$$R_q^2(S_1, S_2) = \frac{1}{q} \text{tr}(P_{S_1} P_{S_2}) \quad \text{Squared Trace Correlation, 1 for } q=p$$

$$R_{o,q}^2(S_1, S_2) = R_q^2(S_1^\perp, S_2^\perp) \quad \text{complement version, 1 for } q=0$$

Bootstrap Scheme: For each “working dimension” $1 \leq q \leq p - 1$:

- Estimate \hat{S}_q
- For $j = 1, \dots, B$ bootstrap replicates, estimate $\hat{S}_q^{(j)}$ and compute $R_q^{2(j)} = R_q^2(\hat{S}_q, \hat{S}_q^{(j)})$, $R_{o,q}^{2(j)} = R_{o,q}^2(\hat{S}_q, \hat{S}_q^{(j)})$, and $\bar{R}_q^{2(j)} \bar{R}_{o,q}^{2(j)}$.
- Calculate the averages \bar{R}_q^2 , $\bar{R}_{o,q}^2$, and $\overline{\bar{R}_q^2 \bar{R}_{o,q}^2}$.

\bar{R}_q^2 , $\bar{R}_{o,q}^2$ and $\overline{\bar{R}_q^2 \bar{R}_{o,q}^2}$ all measure ‘stability’ in estimating S_q .

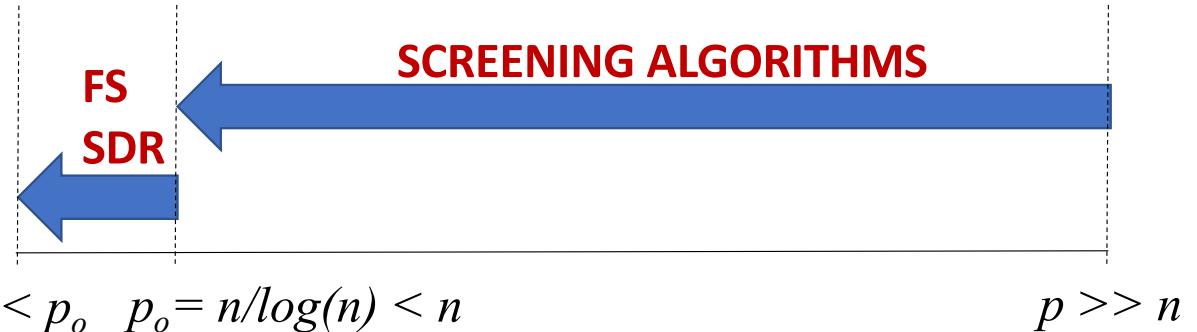
Select \hat{q} where $\overline{\bar{R}_q^2 \bar{R}_{o,q}^2}$ has a “peak”.

What to do when the feature space becomes very high dimensional?

Cut p based on *marginal utilities*, e.g.

$$\omega_j = \text{corr}(Y, X_j) \quad j = 1, 2 \dots p$$

Disregard associations among the X 's but with $p \gg n$ our ability to account for them is too poor!



Covariate Information Number (one feature at a time)

$$\omega_j = \mathbb{C}_{X_j} = \int \mathbb{F}_{x_j} f(x_j) dx_j = \mathbb{J}_{X_j|Y} - \mathbb{J}_{X_j} \quad j = 1, 2 \dots p$$

$$\int \left[\frac{\partial}{\partial x_j} \log f(y | x_j) \right]^2 f(y | x_j) dy$$



Theory and Methods
Covariate Information Number for Feature Screening in Ultrahigh-Dimensional Supervised Problems

Debmalya Nandy, Francesca Chiaromonte & Runze Li

Received 01 Nov 2018, Accepted 08 Dec 2020, Accepted author version posted online: 16 Dec 2020, Published online: 10 Feb 2021

[Download citation](#) <https://doi.org/10.1080/01621459.2020.1864380>



Covariate Information Number for Feature Screening in Ultrahigh-Dimensional Supervised Problems

Abstract: Contemporary high-throughput experimental and surveying techniques give rise to ultrahigh-dimensional supervised problems with sparse signals; that is, a limited number of observations (n), each with a very large number of covariates ($p \gg n$), only a small share of which is truly associated with the response. In these settings, major concerns on computational burden, algorithmic stability, and statistical accuracy call for substantially reducing the feature space by **eliminating redundant covariates before the use of any sophisticated statistical analysis**. Along the lines of *Pearson's correlation coefficient-based sure independence screening* and other model- and correlation-based feature screening methods, we propose a model-free procedure called *covariate information number-sure independence screening* (CIS). CIS uses a **marginal utility connected to the notion of the traditional Fisher information**, possesses the **sure screening property**, and is applicable to **any type of response (features) with continuous features (response)**. Simulations and an application to transcriptomic data on rats reveal the comparative strengths of CIS over some popular feature screening methods.