

NONPARAMETRIC SCREENING UNDER CONDITIONAL STRICTLY CONVEX LOSS FOR ULTRAHIGH DIMENSIONAL SPARSE DATA

BY XU HAN

Temple University

Sure screening technique has been considered as a powerful tool to handle the ultrahigh dimensional variable selection problems, where the dimensionality p and the sample size n can satisfy the NP dimensionality $\log p = O(n^a)$ for some $a > 0$ [*J. R. Stat. Soc. Ser. B. Stat. Methodol.* **70** (2008) 849–911]. The current paper aims to simultaneously tackle the “universality” and “effectiveness” of sure screening procedures. For the “universality,” we develop a general and unified framework for nonparametric screening methods from a loss function perspective. Consider a loss function to measure the divergence of the response variable and the underlying nonparametric function of covariates. We newly propose a class of loss functions called conditional strictly convex loss, which contains, but is not limited to, negative log likelihood loss from one-parameter exponential families, exponential loss for binary classification and quantile regression loss. The sure screening property and model selection size control will be established within this class of loss functions. For the “effectiveness,” we focus on a goodness-of-fit nonparametric screening (Goffins) method under conditional strictly convex loss. Interestingly, we can achieve a better convergence probability of containing the true model compared with related literature. The superior performance of our proposed method has been further demonstrated by extensive simulation studies and some real scientific data example.

1. Introduction. Ultrahigh-dimensional variable selection has become an important problem in modern statistical research due to the big data collection in a variety of scientific areas, such as genomics, bioinformatics, functional magnetic resonance imaging, high frequency finance, etc. In all of these problems, statisticians want to select the important covariates associated with the response variable from p covariates. However, the dimensionality p can grow much faster than the sample size n . More specifically, $\log p = O(n^a)$ for some $a > 0$, which is denoted as nonpolynomial order (NP) [Fan and Lv (2008)]. As Fan, Samworth and Wu (2009) has pointed out: existing variable selection methods based on penalized

Received June 2017; revised June 2018.

¹Supported in part by the summer research funding from Fox Business School at Temple University.

MSC2010 subject classifications. 62G99.

Key words and phrases. Ultrahigh dimensional variable selection, sure screening property, goodness-of-fit nonparametric screening, conditional strictly convex loss.

pseudo likelihood estimation [e.g., Tibshirani (1996), Fan and Li (2001), Zou and Hastie (2005), Zou (2006), Candès and Tao (2007), Zou and Li (2008), Zhang (2010)] can suffer from the simultaneous challenges to computational expediency, statistical accuracy and algorithmic stability in ultrahigh dimensional problems.

To handle the challenges in the ultrahigh-dimensional problems, Fan and Lv (2008) introduced a new statistical framework, sure independence screening. Their original method focused on the Gaussian linear regression models, and the important predictors were selected via the marginal correlation ranking. Formally, let M_\star be the set of true important variables, and \widehat{M}_n be the selected variables based on some procedure, then

$$(1) \quad P(M_\star \subset \widehat{M}_n) \geq 1 - \varepsilon_n,$$

where $\varepsilon_n > 0$ and $\varepsilon_n \rightarrow 0$ as $n \rightarrow \infty$. This is called the “sure screening property.” Furthermore, the model selection size can be controlled at a polynomial rate of sample size with probability approaching 1. Because of its powerful performance and computational convenience in ultrahigh dimensional problems, the sure screening framework has received increasing attention in the past few years. Existing literature in this framework have mainly focused on the “universality” of the screening procedures, that is, developing procedures for various scenarios which possess the “sure screening property,” for example, generalized linear model by Fan and Song (2010), nonparametric additive model by Fan, Feng and Song (2011), rank based model-free feature screening by Zhu et al. (2011), Cox model by Zhao and Li (2012), robust rank correlation screening by Li et al. (2012), varying coefficient model by Fan, Ma and Dai (2014), empirical likelihood based screening by Chang, Tang and Wu (2013), quantile-adaptive screening by He, Wang and Hong (2013), censored rank independence screening by Song et al. (2014), fused Kolmogorov filter by Mai and Zou (2015). On the other hand, formal pursuit of “effectiveness” of sure screening procedures have been largely ignored. Intuitively, for the convergence probability $1 - \varepsilon_n$ in (1), if ε_n converges to 0 slower, the corresponding screening procedure will have larger possibility of not selecting the true important variables. More specifically, existing literature commonly show that

$$(2) \quad P(M_\star \subset \widehat{M}_n) \geq 1 - s_n \{b \exp(-cn^a)\},$$

where a, b, c are positive values and s_n is the size of true model. The rate a controls how high dimensionality the screening procedure can handle. It will be illustrated in detail in later sections. For a larger a , the probability of containing true model converges to 1 faster. The effect of c can be negligible for a larger a and a sufficient large n . The constant b is not crucial in the asymptotic sense, but it is important for finite sample situations. With the same value of a , a larger value of b indicates that important variables can be mis-selected with higher probability. For some existing results, b can even grow as n increases. Therefore, the constant b and the convergence rate a can be viewed as a measure of effectiveness of a screening method.

Correspondingly, a sure screening procedure will be considered more effective with a larger a and a smaller b in the convergence probability of (2). Although the existing screening procedures have been proved to possess the “sure screening property,” the established convergence of containing the true model can be slow subject to various specific model settings and conditions.

Our first goal in the current paper is to develop a general and unified framework for sure screening methods from a loss function perspective. Consider a response variable Y which distribution depends on parameter θ . Suppose θ is a function of p -dimensional covariate vector $\mathbf{X} = (X_1, \dots, X_p)^T$. We are interested in selecting the covariates X_j 's which are associated with the response variable Y through a nonparametric function $\theta = f(X_1, \dots, X_p)$. For notational convenience, we will write as $\theta(\mathbf{X})$ to denote its dependence on the covariates \mathbf{X} . In later presentation, we sometimes simply write it as θ for the true function, and the readers should be reminded that the θ is a function on \mathbf{X} . This setting includes a variety of commonly used regression models:

EXAMPLE 1 (Gaussian regression). Assume that $Y|\mathbf{X} = \mathbf{x}$ is from $N(\theta(\mathbf{x}), \sigma^2)$ for some constant $\sigma > 0$.

EXAMPLE 2 (Logistic regression). Assume that $Y|\mathbf{X} = \mathbf{x}$ is from Bernoulli distribution and $\ln P(Y = 1|\mathbf{X} = \mathbf{x}) - \ln P(Y = 0|\mathbf{X} = \mathbf{x}) = \theta(\mathbf{x})$.

EXAMPLE 3 (Poisson regression). Assume that $Y|\mathbf{X} = \mathbf{x}$ is from Poisson distribution and $\ln E(Y|\mathbf{X} = \mathbf{x}) = \theta(\mathbf{x})$.

EXAMPLE 4 (Quantile regression). Let $Q_\alpha(Y|\mathbf{X} = \mathbf{x})$ be the α th quantile of the distribution for $Y|\mathbf{X} = \mathbf{x}$, then assume $Q_\alpha(Y|\mathbf{X} = \mathbf{x}) = \theta(\mathbf{x})$.

The above Examples 1–3 fall within the general framework of mean regression:

$$(3) \quad E(Y|\mathbf{X} = \mathbf{x}) = h(\theta(\mathbf{x})) = g^{-1}(f(x_1, \dots, x_p)),$$

where h is some known function, f is a nonparametric function and g is called the link function. When g is the canonical link, that is, $g = (h)^{-1}$, we have $\theta(\mathbf{X}) = f(X_1, \dots, X_p)$. However, Example 4 is different from the mean regression.

The above regression models are equivalent to considering a loss function $l(\omega, Y)$ for measuring the divergence between a generic variable ω and the response variable Y where ω is a function of \mathbf{X} , and assuming that the true model of θ will minimize $E[l(\omega, Y)|\mathbf{X} = \mathbf{x}]$ with respect to ω . For instance, in the above Examples 1–3, we can choose $l(\omega, Y)$ as the negative of the log-likelihood of $Y|\mathbf{X} = \mathbf{x}$; In the above Example 4, we can choose $l(\omega, Y) = (Y - \omega)[\alpha - \mathbf{I}(Y - \omega < 0)]$, where \mathbf{I} is an indicator function. Therefore, we will select the important covariates X_j 's associated with Y based on such a loss function.

In the current paper, we newly propose a definition of loss function called conditional strictly convex loss, which contains, but is not limited to, negative log-likelihood loss for one-parameter exponential families, exponential loss for binary classification and quantile regression loss for robust estimation. Our sure screening property is established within such a wide class of loss functions. Therefore, several existing screening methods automatically fall within our framework, including Fan, Feng and Song (2011) for nonparametric additive models and He, Wang and Hong (2013) for quantile regression, although their proposed screening procedures can be different from ours. In addition, many more screening methods are suggested by our framework, for example, generalized additive models, binary classification by exponential loss and so on.

Our second goal of the current paper is to develop screening methods under conditional strictly convex loss with better convergence probability of containing the true model. We treat the marginal regression as fitting the response variable with componentwise covariates via the loss function. We impose an additive model structure for the unknown nonparametric function approximated by B-spline basis. Interestingly, if we consider the goodness-of-fit statistics as the marginal utility to rank the importance of each covariate to the joint model, we can achieve a much better convergence probability of containing the true model compared with other related literature. Detailed comparison between our results with other related literature will be presented in Section 3. Furthermore, our selected model size can be controlled at the level of sample size n rather than the dimensionality p_n with high probability.

The major contribution of the current paper is to simultaneously tackle the issues of “universality” and “effectiveness.” For the “universality,” we establish the sure screening property within a unified framework through the introduction of a new class of loss functions: conditional strictly convex loss; For the “effectiveness,” within this framework, we show that the goodness-of-fit nonparametric screening methods can achieve a better convergence probability of containing the true model compared with related literature.

Theoretical pursuit of “universality” and “effectiveness” for screening procedures in the current paper has shed new light on the choice of sure screening methods and greatly benefited the applications of screening methods in practice. For example, the superior performance of our proposed method compared with other existing screening procedures will be further demonstrated by extensive simulation studies and some real scientific data example. Our method is called goodness-of-fit nonparametric screening (Goffins). To stabilize the computation performance, we also provide an iterative screening procedure and an improved variant to handle the situations where covariates are possibly correlated.

The rest of this paper will be organized as follows: Section 2 introduces the conditional strictly convex loss, the B-spline approximation and the goodness of fit nonparametric screening; Section 3 establishes the exponential bound, the sure screening properties and the control of model selection size; Section 4 proposes

an iterative screening procedure and an improved variant; Section 5 provides simulation studies and real data analysis. All the technical proofs and some numerical results are relegated to the Supplementary Material [Han (2019)].

2. Nonparametric screening under convex loss.

2.1. *Conditional strictly convex loss.* Let $l(x, y) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be a function and assume the partial derivative $\partial l(x, y)/\partial x$ exists almost everywhere for x throughout the paper. We consider $l(\omega, Y)$ as a loss function to measure the divergence between a generic variable ω and the response variable Y . We assume the convexity of $l(\omega, Y)$ in the ω position, that is, $l(t_1\omega_1 + t_2\omega_2, Y) \geq t_1l(\omega_1, Y) + t_2l(\omega_2, Y)$ for any real values $t_1 + t_2 = 1$ and $t_1, t_2 > 0$. Here, ω is a function of covariates \mathbf{X} , and can be written as $\omega(\mathbf{X})$ to denote its dependence on \mathbf{X} . For notational convenience, we sometimes simply write it as ω . Suppose the distribution of Y depends on some parameter θ where θ is a nonparametric function of the covariates \mathbf{X} . Assume the true model of θ minimizes $E[l(\omega, Y)|\mathbf{X}]$ with respect to ω .

In the current paper, we will newly propose a definition of loss function called conditional strictly convex loss. Our sure screening method will be established within such a wide class of convex loss functions.

DEFINITION 1. If $\partial E[l(\omega, Y)|\mathbf{X}]/\partial \omega$ is continuously differentiable in ω and $\partial^2 E[l(\omega, Y)|\mathbf{X}]/\partial \omega^2 > 0$, then $l(\omega, Y)$ is called a conditional strictly convex loss function.

The conditional strictly convex loss includes, but is not limited to, the following three major types of loss functions:

Type 1: Negative log-likelihood loss for exponential families. Suppose that the random variable Y is from a one-parameter exponential family with density function

$$(4) \quad f_{Y|X}(y; \theta) = \exp(y\theta - b(\theta) + c(y))$$

for some known functions $b(\cdot)$ and $c(\cdot)$ where $b''(\cdot)$ exists. Consider the negative log-likelihood loss:

$$(5) \quad l(\omega, Y) = -[\omega Y - b(\omega) + c(Y)].$$

Minimization of $E[l(\omega, Y)|\mathbf{X}]$ with respect to ω and letting θ be the minimizer leads to $E[Y|\mathbf{X}] = b'(\theta)$, which naturally belongs to the mean regression (3). This is the setting of generalized additive model in Stone (1986). Note that the second derivative of $l(\omega, Y)$ with respect to ω is $b''(\omega)$, and $b''(\theta)$ is the variance of Y from the exponential families.

The loss function (5) can be better understood by some popular regression models:

EXAMPLE 1 (Gaussian regression). $b(\theta) = \theta^2/2$, $c(y) = -y^2/2$ and $l(\omega, Y) = (Y - \omega)^2/2$.

EXAMPLE 2 (Logistic regression). $b(\theta) = \ln(1 + \exp(\theta))$, $c(y) = 0$ and $l(\omega, Y) = -\omega Y + \ln(1 + \exp(\omega))$.

EXAMPLE 3 (Poisson regression). $b(\theta) = \exp(\theta)$, $c(y) = -\ln(y!)$ and $l(\omega, Y) = -Y\omega + \exp(\omega) + \ln(Y!)$.

Type 2: Exponential loss for classification. In classification problems, suppose $Y \in \{-1, 1\}$ and $P(Y = 1|\mathbf{X} = \mathbf{x}) = p(\mathbf{x})$. The goal is to construct a classifier $\theta(\mathbf{x})$. When new covariates \mathbf{X} are available, predict the corresponding class type Y as 1 if $\theta(\mathbf{X}) > c$ and as -1 if $\theta(\mathbf{X}) < c$ where c is some threshold. The exponential loss is defined as

$$(6) \quad l(\omega, Y) = \exp(-Y\omega),$$

which has been considered as a smooth approximation to the misclassification loss [Freund and Schapire (1997)]. Minimization of $E[l(\omega, Y)|\mathbf{X}]$ with respect to ω and letting θ be the minimizer leads to

$$\ln \frac{P(Y = 1|\mathbf{X})}{P(Y = -1|\mathbf{X})} = 2\theta.$$

Type 3: Quantile regression loss. For many practical problems, the distribution information of response variable Y is usually not available or complicated. Instead of imposing a full distribution, quantile regression framework assumes that the α th quantile of Y given \mathbf{X} , $Q_\alpha(Y|\mathbf{X})$, is some function of \mathbf{X} , thus the distribution assumption can be substantially relaxed [Koenker (2005)]. Correspondingly, consider the loss function

$$(7) \quad l(\omega, Y) = (Y - \omega)\{\alpha - \mathbf{I}(Y - \omega < 0)\}$$

for $0 < \alpha < 1$ where \mathbf{I} is an indicator function. When $\alpha = 1/2$, this is proportional to the least absolute deviation loss $|Y - \omega|$, which is popularly used for robust regression. The loss function $l(\omega, Y)$ is not differentiable in ω . This is a key difference from the aforementioned loss functions. Minimization of $E[l(\omega, Y)|\mathbf{X}]$ with respect to ω yields $Q_\alpha(Y|\mathbf{X}) = \theta$ where θ is the minimizer.

The following Proposition 2.1 shows that with mild conditions, Types 1–3 belong to the conditional strictly convex loss.

PROPOSITION 2.1. *For Type 1, if b'' is strictly positive and is a continuous function, then (5) belongs to the conditional strictly convex loss; for Type 2, (6) belongs to the conditional strictly convex loss; for Type 3, if the conditional distribution of $Y|\mathbf{X}$ has a continuous density function $f_{Y|\mathbf{X}}$ and $f_{Y|\mathbf{X}} > 0$ on any bounded domain, then (7) belongs to the conditional strictly convex loss.*

Without any further investigation, one might simply group Types 1 and 2 in Proposition 2.1 as one class since the corresponding loss functions are second differentiable in ω . However, we will show in Section 3.4 that even for Types 1 and 2, the loss functions possess some fundamental differences in the underlying structures, which raises challenges for proving the model selection size control in Section 3.4.

The name of conditional strictly convex loss is borrowed from “strictly convex function.” However, there are some major differences between the two concepts. If $l(x, y)$ is a strictly convex function in x and $l'(x, y)$ is continuously differentiable in x , then l is also a conditional strictly convex loss, but a conditional strictly convex loss might not be a strictly convex function; see Type 3 quantile regression loss as such a counterexample.

A class of convex loss, Bregman divergence, can also be considered here. For a given convex function $q(\cdot)$ with derivative $q'(\cdot)$, the Bregman divergence [Brègman (1967)] is defined as

$$(8) \quad l(\omega, Y) = q(\omega) - q(Y) + (Y - \omega)q'(Y).$$

Note that $l(\omega, Y)$ is not generally a symmetric function in ω and Y . Suppose $q'(\cdot)$ is continuously differentiable and $q''(\cdot) > 0$, it is easy to show that such Bregman divergence belongs to the conditional strictly convex loss. It is impossible for us to list all the possibilities here, thus we will not go any further in this direction. It is worth mentioning that the quantile regression loss (7) does not belong to Bregman divergence. More detailed discussions about Bregman divergence are referred to Zhang, Jiang and Shang (2009).

2.2. Goodness-of-fit nonparametric screening. To capture the nonparametric structure of $\theta(\mathbf{X})$, a powerful model for dimensionality reduction is the additive model:

$$(9) \quad \theta(\mathbf{X}) = m_1(X_1) + \cdots + m_p(X_p) + \mu,$$

where $m_j(\cdot)$ are the square integrable functions and μ is an unknown constant. For identifiability, we assume $E[m_j(X_j)] = 0$ for $j = 1, \dots, p$. Let $M_\star = \{j : E[m_j(X_j)]^2 > 0\}$ be the true sparse model with nonsparsity size $s_n = |M_\star|$. Suppose we have observed data $\{(\mathbf{X}_i, Y_i)\}$ for $i = 1, \dots, n$, which are independent copies of (\mathbf{X}, Y) . The dimensionality p is ultrahigh and satisfies $\log p = O(n^a)$ for some $a > 0$. Based on the sample data, we aim to select a subset of covariates which contains M_\star with moderate size. We allow p to grow with n , and denote the dimensionality as p_n .

In this paper, we refer to marginal regression as fitting models with componentwise covariates through the loss function $l(\omega, Y)$. We define the population version of the minimizer of the componentwise regression as

$$(10) \quad f_j^M(X_j) \equiv \arg \min_{f_j \in L_2(P)} E[l(f_j(X_j), Y)],$$

where P denotes the joint distribution of (\mathbf{X}, Y) and $L_2(P)$ is the class of square integrable functions under measure P . We use B-spline basis to approximate the marginal nonparametric regression function. Let S_n be the space of polynomial splines of degree $l \geq 1$. Stone (1986) has shown that under some smoothness conditions, the nonparametric functions can be well approximated by functions in S_n . Correspondingly, we define

$$(11) \quad f_{nj}^M(X_j) \equiv \arg \min_{f_j \in S_n} E[l(f_j(X_j), Y)].$$

We also define the marginal minimum divergence estimator as

$$(12) \quad \widehat{f}_{nj}^M(X_j) \equiv \arg \min_{f_j \in S_n} \mathbb{P}_n l(f_j(X_j), Y),$$

where $\mathbb{P}_n g(\mathbf{X}, Y) = n^{-1} \sum_{i=1}^n g(\mathbf{X}_i, Y_i)$ is the empirical expectation for generic function $g(\cdot)$. Let $\{\Psi_k\}_{k=1}^{d_n}$ denote a normalized B-spline basis with $\|\Psi_k\|_\infty \leq 1$, where $\|\cdot\|_\infty$ is the sup norm. For any $f_{nj} \in S_n$, we have

$$(13) \quad f_{nj}(x) = \sum_{k=1}^{d_n} \Psi_k(x) \beta_{jk}, \quad 1 \leq j \leq p$$

for some coefficients $\{\beta_{jk}\}_{k=1}^{d_n}$. The construction of the B spline basis can be found in the well-known books, for example, de Boor (1978). Let $\Psi_j \equiv \Psi_j(X_j) = (\Psi_1(X_j), \dots, \Psi_{d_n}(X_j))^T$, therefore, we can express

$$(14) \quad f_{nj}^M(X_j) = \Psi_j^T \beta_j^M, \quad \widehat{f}_{nj}^M(X_j) = \Psi_j^T \widehat{\beta}_j^M,$$

where β_j^M and $\widehat{\beta}_j^M$ are the d_n dimensional coefficient vector for the minimizers of (11) and (12).

We will consider a sure screening procedure based on goodness-of-fit statistics. Formally, let

$$G_{n,j} = \mathbb{P}_n \{l(\widehat{\beta}_0^M, Y) - l(\Psi_j^T \widehat{\beta}_j^M, Y)\}, \quad j = 1, \dots, p_n,$$

where $\widehat{\beta}_0^M \equiv \arg \min_{\beta_0 \in \mathbb{R}} \mathbb{P}_n l(\beta_0, Y)$. Correspondingly, for the population level,

$$G_j^* = E \{l(\beta_0^M, Y) - l(\Psi_j^T \beta_j^M, Y)\}, \quad j = 1, \dots, p_n,$$

where $\beta_0^M \equiv \arg \min_{\beta_0 \in \mathbb{R}} E l(\beta_0, Y)$. The goodness-of-fit statistics compare the marginal regression model with the null model (no variables included in the model). Intuitively, if the marginal contribution of an individual variable is significant to the response variable, the goodness-of-fit measure should be relatively large. We select model by $\widehat{M}_{v_n} = \{1 \leq j \leq p_n : G_{n,j} \geq v_n\}$ for a predetermined threshold v_n . Our screening method is called goodness-of-fit nonparametric screening (Goffins). We intentionally use the letter ‘‘G’’ in $G_{n,j}$ and G_j^* to denote the goodness-of-fit statistics.

When l is the squared error loss, since the term $\mathbb{P}_n l(\widehat{\beta}_0^M, Y)$ in $G_{n,j}$ is not affected by the index j , Goffins is equivalent to screening based on the sum of squared residuals, that is, select the model by $\{1 \leq j \leq p_n : \mathbb{P}_n(Y - \Psi_j^T \widehat{\beta}_{j1}^M)^2 \leq \mu_n\}$ for some threshold $\mu_n > 0$. Note that, $\mathbb{P}_n(Y - \Psi_j^T \widehat{\beta}_j^M)^2$ can be further expressed as $\mathbb{P}_n Y^2 - \mathbb{P}_n(\Psi_j^T \widehat{\beta}_j^M)^2$. Therefore, Goffins under the squared error loss is equivalent to selecting the model by $\{1 \leq j \leq p_n : \mathbb{P}_n(\Psi_j^T \widehat{\beta}_j^M)^2 \geq \gamma_n\}$ for some threshold $\gamma_n > 0$. More generally, when l is the negative log-likelihood loss for exponential families, Goffins is equivalent to screening based on the likelihood ratio statistics. For parametric model based likelihood ratio screening; see Fan and Song (2010).

3. Sure screening properties.

3.1. *Preliminaries.* In this paper, we will show that our goodness-of-fit nonparametric screening (Goffins) has the sure screening property, and the number of the selected variables has moderate size. Let $[a, b]$ be the support of covariates X_j . The following conditions are needed:

(A) The nonparametric marginal functions $\{f_j^M\}_{j=1}^p$ belong to a class of functions \mathfrak{F} whose r th derivative $f^{(r)}$ exists and is Lipschitz of order α :

$$(15) \quad \mathfrak{F} = \{f(\cdot) : |f^{(r)}(s) - f^{(r)}(t)| \leq K|s - t|^\alpha, \text{ for } s, t \in [a, b]\}$$

for some positive constant K , where r is a nonnegative integer and $\alpha \in (0, 1]$ such that $d = r + \alpha > 0.5$.

(B) The marginal density functions g_j of X_j satisfies $0 < K_1 \leq g_j(X_j) \leq K_2 < \infty$ on $[a, b]$ for $1 \leq j \leq p$ for some constants K_1 and K_2 .

(C) The unknown nonparametric function $\theta(\mathbf{X})$ satisfies that $\sup_{\mathbf{X} \in \mathbb{R}^{p_n}} |\theta(\mathbf{X})| < M$ from some positive constant M .

Conditions A, B and C are standard regularity assumptions for nonparametric regression in Stone (1986), Fan, Feng and Song (2011), He, Wang and Hong (2013), etc.

The following Lemma 3.1 shows that the approximation error of marginal regression f_{nj}^M in (11) to marginal nonparametric projection f_j^M in (10) is negligible.

LEMMA 3.1. *If l is a conditional strictly convex loss, under Conditions A–C, assume that f_j^M is uniformly bounded for $j = 1, \dots, p$, then there exists a positive constant C_1 such that $E(f_j^M - f_{nj}^M)^2 \leq C_1 d_n^{-2d}$, where d is defined in Condition A.*

To show that for $j \in M_\star$, G_j^\star has a nonvanishing signal, we also need the following conditions:

(D) $\min_{j \in M_\star} E[f_j^M(X_j) - Ef_j^M(X_j)]^2 \geq c_1 d_n n^{-2\kappa}$, for some $0 < \kappa < d/(2d + 1)$ and $c_1 > 0$.

(E) $d_n^{-2} \leq c_1(1 - \xi)^2 n^{-2\kappa} / 4C_1$ for some $\xi \in (0, 1)$.

Condition D requires that the marginal nonparametric projections are at a certain strength level separate from the noise. Therefore, we can select the significant covariates based on a threshold. Similar conditions also appear in related literature on nonparametric screening, for example, [Fan, Feng and Song \(2011\)](#) and [He, Wang and Hong \(2013\)](#). See detailed discussion in Section 2 of Supplementary Material [[Han \(2019\)](#)].

LEMMA 3.2. *Under conditions in Lemma 3.1, in addition, Condition D and E are also satisfied, then $\min_{j \in M_\star} G_j^\star \geq \frac{b^\star}{2} c_1 \xi d_n n^{-2\kappa}$ for some positive constant b^\star .*

As we will show in later sections, the sure screening property depends on the characteristics of a generalized definition for partial derivative of loss function $l(x, y) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ with respect to x . More specifically, let $\tilde{l}(x, y) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be a Riemann integrable function with respect to x such that for any $x_1 > x_2$ and any y , $l(x_1, y) - l(x_2, y) = \int_{x_2}^{x_1} \tilde{l}(s, y) ds$. Since $l(x, y)$ is differentiable in x almost everywhere, such \tilde{l} exists and is unique almost everywhere in x . For notational convenience, we simply use $l'(x, y)$ to denote one such $\tilde{l}(x, y)$. When $l(x, y)$ is differentiable in x , $l'(x, y)$ is uniquely determined. When we consider the quantile regression loss $l(x, y) = (y - x)\{\alpha - \mathbf{I}(y - x < 0)\}$, if $x > y$, then $l(x, y) = (y - x)(\alpha - 1)$; if $x < y$, then $l(x, y) = (y - x)\alpha$. Except at $x = y$, $\partial l(x, y) / \partial x = \mathbf{I}(y - x < 0) - \alpha$. Hence, for quantile regression loss, for any $x_1 > x_2$ and any y , we have $l(x_1, y) - l(x_2, y) = \int_{x_2}^{x_1} [\mathbf{I}(y - s < 0) - \alpha] ds$. Therefore, we will use $l'(x, y) = \mathbf{I}(y - x < 0) - \alpha$ for the quantile regression loss throughout the paper. The above argument motivates the following Definition 2.

DEFINITION 2. The notation $l'(\omega, Y)$ is defined as follows: for Type 1 and 2, $l'(\omega, Y) = \frac{\partial l(\omega, Y)}{\partial \omega}$; for Type 3, $l'(\omega, Y) = \mathbf{I}(Y - \omega < 0) - \alpha$.

To simplify the discussion, for the loss function $l(x, y)$ which is not differentiable in x but is differentiable in x almost everywhere, we only focus on the quantile regression loss here. However, similar argument also applies to other loss functions beyond quantile regression loss.

To characterize $l'(\omega, Y)$ for the exponential tail bound in Section 3.2, we also need the following definition for sub-Gaussian random variables.

DEFINITION 3. A random variable X with mean $\mu = EX$ is called σ -sub-Gaussian if there is a positive number σ such that

$$E \exp(\lambda(X - \mu)) \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right) \quad \forall \lambda \in \mathbb{R}.$$

Note that if $X \sim N(\mu, \sigma^2)$, then X is σ -sub-Gaussian. If random variable X is bounded such that $a \leq X \leq b$, then X is sub-Gaussian with $\sigma = (b - a)/2$; see Buldygin and Kozachenko (2000) for more details.

3.2. *Exponential bound for marginal minimum divergence estimator.* Since both \widehat{f}_{nj}^M and f_{nj}^M can be expressed in terms of B-spline basis functions, it is crucial to establish an exponential bound for the tail probability of $\|\widehat{\beta}_j^M - \beta_j^M\|$. The sharpness of this exponential bound directly affect the convergence probability of the screening method.

The following Theorem 3.1 provides an exponential bound for the tail probability of marginal minimum divergence estimator for the B-spline coefficients. It will serve as the cornerstone for our later derivations of the other theorems. The following conditions are required for Theorem 3.1:

- (F) $d_n = o(n^{1/3})$ and $d_n = O(n^{2\kappa})$.
- (G) $E[l'(\omega, Y)|\mathbf{X}]$ is bounded for any bounded ω .

PROPOSITION 3.1. *For Types 1–3, under the conditions in Proposition 2.1, condition G is satisfied.*

The tail probability of $\|\widehat{\beta}_j^M - \beta_j^M\|$ depends on the properties of $l'(\omega, Y)$. More specifically, we will consider the following set of conditions:

- (H1) $l'(\omega, Y)$ is bounded for any bounded ω ;
- (H2) $l'(\omega, Y)$ conditional on \mathbf{X} is a σ -sub-Gaussian random variable where σ does not depend on \mathbf{X} ;
- (H3) For any bounded ω , $E[\exp(\lambda l'(\omega, Y))|\mathbf{X}] < \infty$ for all $|\lambda| \leq c_0$ with some constant $c_0 > 0$.

The notation $l'(\omega, Y)$ in Condition G and H1–H3 is based on Definition 2. By Definition 3, if Condition H1 is satisfied, then Condition H2 is also satisfied; if Condition H2 is satisfied, then Condition H3 is also satisfied. In the following Theorem 3.1, we will show that with stronger assumption a better tail probability can be correspondingly achieved.

To better understand the wide applicability of Conditions H1–H3, let us consider some examples from Types 1–3 which satisfy these conditions. Some popular regression models can be summarized in the following Proposition 3.2.

PROPOSITION 3.2. *Types 2 and 3 satisfy Condition H1. For Type 1, if $Y|\mathbf{X}$ follows Bernoulli distribution, then (5) satisfies Condition H1; if $Y|\mathbf{X}$ follows Normal distribution, then (5) satisfies Condition H2; if $Y|\mathbf{X}$ follows Poisson distribution, then (5) satisfies Condition H3.*

Furthermore, if $Y|\mathbf{X}$ follows some other distributions in the exponential family, under some regularity conditions, it is possible that the corresponding loss function (5) also satisfy Condition H3. For example, if $Y|\mathbf{X} \sim \text{Laplace}(\mu(\mathbf{X}), b)$ with a known parameter b , then Condition H3 is satisfied. If $Y|\mathbf{X} \sim \text{Exponential}(\lambda(\mathbf{X}))$, and if there exists a positive constant c such that $\lambda(\mathbf{X}) \geq c$, then Condition H3 is satisfied. Similar arguments for verifying Condition H3 also apply to Chi-square distribution, negative binomial distribution, inverse-Gaussian distribution with a known shape parameter, Gamma distribution with a known scale parameter. To save space, we will not discuss in detail for these examples.

THEOREM 3.1. *For a convex loss $l(\omega, Y)$, if it is also a conditional strictly convex loss, for any constant $c_3 > 0$, under Conditions C, F, G, there exist positive constants c_4 and c_5 such that for sufficiently large n :*

if Condition H1 is satisfied and $d_n = o(n^{1-2\kappa})$, then

$$(16) \quad P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_3 d_n n^{-2\kappa}) \leq \exp(-c_4 n^{1-2\kappa} d_n^{-1});$$

if Condition H2 is satisfied and $d_n = o(n^{(1-2\kappa)/2})$, then

$$(17) \quad P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_3 d_n n^{-2\kappa}) \leq \exp(-c_4 n^{1-2\kappa} d_n^{-1}) + \exp(-c_5 n^{1-2\kappa} d_n^{-2});$$

if Condition H3 is satisfied and $d_n = o(n^{(1-2\kappa)/3})$, then

$$(18) \quad P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_3 d_n n^{-2\kappa}) \leq \exp(-c_4 n^{1-2\kappa} d_n^{-1}) + 2 \exp(-c_5 n^{1/2-\kappa} d_n^{-3/2}).$$

In (16), when $d_n = o(n^{1-2\kappa})$, $n^{1-2\kappa} d_n^{-1}$ in the tail probability diverges to infinity as n increases, which implies that the tail probability converges to zero. Similar arguments also apply to (17) and (18). It is worth mentioning that Theorem 3.1 is proved based on a unified argument with some modifications according to each situation of Conditions H1–H3. The proof is different from the related literature and can be of independent research interest.

3.3. Sure screening. Based on Theorem 3.1 for estimation of B-spline coefficients, we are now ready to establish the sure screening property for our Goffins method. Different properties of loss functions can lead to different convergence probabilities of containing the true model.

THEOREM 3.2. *Under the conditions in Theorem 3.1 and Lemma 3.2:*

(i) *for Types 1 and 2, there exists a positive constant ζ , then by taking $v_n = \nu d_n n^{-2\kappa}$ with $0 < \nu \leq \zeta$, there exists positive constants c_4, c_5 and c_6 such that if Condition H1 is satisfied and $d_n = o(n^{1-2\kappa})$, then*

$$(19) \quad P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n [\exp(-c_4 n^{1-2\kappa} d_n^{-1}) + 6 \exp(-c_5 n^{1-2\kappa})];$$

if Condition H2 is satisfied and $d_n = o(n^{(1-2\kappa)/2})$, then

$$(20) \quad P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n [\exp(-c_4 n^{1-2\kappa} d_n^{-1}) + \exp(-c_5 n^{1-2\kappa} d_n^{-2}) + 6 \exp(-c_6 n^{1-2\kappa})];$$

if Condition H3 is satisfied and $d_n = o(n^{(1-2\kappa)/3})$, then

$$(21) \quad P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n [\exp(-c_4 n^{1-2\kappa} d_n^{-1}) + 2 \exp(-c_5 n^{1/2-\kappa} d_n^{-3/2}) + 6 \exp(-c_6 n^{1-2\kappa})];$$

(ii) for Type 3, if $d_n = o(n^{1-2\kappa})$, take $v_n = v d_n n^{-2\kappa}$ with $v \leq b^* c_1 \xi / 4$ where b^* is defined in Lemma 3.2, there exist positive constants c_4 and c_5 such that

$$(22) \quad P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n [\exp(-c_4 n^{1-2\kappa} d_n^{-1}) + 12 \exp(-c_5 n^{1-2\kappa})].$$

Theorem 3.2 shows that our Goffins method corresponding to conditional strictly convex loss possesses the sure screening property. It follows from Theorem 3.2 that in (19) and (22) we can handle the NP-dimensionality: $\log p_n = o(n^{1-2\kappa} d_n^{-1})$. Under this condition, $P(M_\star \subset \widehat{M}_{v_n}) \rightarrow 1$ to achieve the sure screening property. For (20), the NP-dimensionality will be changed to $\log p_n = o(n^{1-2\kappa} d_n^{-2})$ and for (21) we can handle $\log p_n = o(n^{1/2-\kappa} d_n^{-3/2})$.

The proof of Theorem 3.2 in the Supplementary Material [Han (2019)] is not limited to Types 1–3. For example, let Class A be the loss functions such that $l''(\omega, Y) \equiv \partial^2 l(\omega, Y) / \partial \omega^2$ exists, $l''(\omega, Y)$ is continuous in ω , $l''(\omega, Y) > 0$ and $l''(\omega, Y)$ is bounded when ω is bounded, then the results corresponding to Types 1–2 in Theorem 3.2 are also valid for the loss functions in Class A. It is not difficult to verify that Types 1–2 are only special examples in Class A. When $l(\omega, Y)$ is not differentiable in ω , the discussion is more complicated. Let $\tilde{l}(x, y) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be a Riemann integrable function with respect to x such that for any $x_1 > x_2$ and any y , $l(x_1, y) - l(x_2, y) = \int_{x_2}^{x_1} \tilde{l}(s, y) ds$. Since we assume that the loss function $l(x, y)$ is differentiable in x almost everywhere, such \tilde{l} exists and is unique almost everywhere in x . Let Class B be the loss functions such that there exists a corresponding \tilde{l} where $\tilde{l}(\omega, Y)$ is bounded and $\tilde{l}(\omega, Y)$ is nondecreasing in ω , then the results corresponding to Type 3 is also valid for the loss functions in Class B. Our definition of $l'(\omega, Y)$ for quantile regression loss clearly satisfies the conditions in Class B.

3.4. *Controlling selection size.* The sure screening methods will not be informative unless the model selection size can be controlled at a reasonable level. The following Theorem 3.3 shows that our Goffins method can control the size of the selected variables at the level of the sample size n rather than the dimension p_n . However, for controlling the model size, our definition of conditional strictly convex loss is not sufficient for the discussion. We need a finer class of loss functions which possesses certain structures. This is the motivation of our following Definition 4.

DEFINITION 4. If a convex loss $l(\omega, Y)$ satisfies that

$$\partial E[l(\omega, Y)|\mathbf{X}]/\partial\omega = G(\omega) - H(\mathbf{X})K(\omega)$$

for some functions $G(\cdot)$, $H(\cdot)$ and $K(\cdot)$, then $l(\omega, Y)$ is called conditional derivative separable loss.

PROPOSITION 3.3. *Types 1–3 are conditional derivative separable loss functions.*

For Type 1, $\partial E[l(\omega, Y)|\mathbf{X}]/\partial\omega = b'(\omega) - \mu(\mathbf{X})$ where $\mu(\mathbf{X}) = E(Y|\mathbf{X})$. Therefore, $G(\omega) = b'(\omega)$, $H(\mathbf{X}) = \mu(\mathbf{X})$ and $K(\omega) = 1$.

For Type 2, $\partial E[l(\omega, Y)|\mathbf{X}]/\partial\omega = \exp(\omega) - [\exp(\omega) + \exp(-\omega)]p(\mathbf{X})$ where $p(\mathbf{X}) = E(Y|\mathbf{X})$. Therefore, $G(\omega) = \exp(\omega)$, $H(\mathbf{X}) = p(\mathbf{X})$ and $K(\omega) = \exp(\omega) + \exp(-\omega)$.

For Type 3, $\partial E[l(\omega, Y)|\mathbf{X}]/\partial\omega = F_{Y|\mathbf{X}}(\omega) - \alpha$ where $F_{Y|\mathbf{X}}$ is the conditional cumulative distribution function of $Y|\mathbf{X}$. Therefore, $G(\omega) = F_{Y|\mathbf{X}}(\omega)$, $H(\mathbf{X}) = \alpha$ and $K(\omega) = 1$.

Detailed discussions reveal different structures of loss functions. For example, Types 1 and 3 both have a constant $K(\omega) = 1$ while Type 2 does not have such property. Furthermore, Types 1 and 2 are second differentiable in ω but Type 3 is not differentiable in ω . Fortunately, we can propose a unified proof to bound $\sum_{j=1}^{p_n} E(f_{nj}^M - Ef_{nj}^M)^2$ when the loss function $l(\omega, Y)$ is conditional derivative separable loss and conditional strictly convex loss, which will serve as a major step for controlling the model selection size.

THEOREM 3.3. *Let $\Psi = (\Psi_1, \dots, \Psi_{p_n})^T$, β^* be the coefficient vector of basis functions for the joint regression model of $\theta(\mathbf{X})$ on \mathbf{X} , β_0^* be the intercept term in the joint regression model and $\Sigma = E\Psi\Psi^T$. If l is a conditional derivative separable loss and conditional strictly convex loss, under conditions in Theorem 3.2, in addition, $E(\Psi^T\beta^*)^2 = O(1)$ and $K(\beta_0^*) \neq 0$, then we have:*

(i) $\sum_{j=1}^{p_n} E(f_{nj}^M - Ef_{nj}^M)^2 = O(d_n\lambda_{\max}(\Sigma))$;

(ii) *with v_n described in Theorem 3.2, there exist constants c_4, c_5, c_6 such that Types 1 and 2: if Condition H1 is satisfied and $d_n = o(n^{1-2\kappa})$, then*

$$P(|\widehat{M}_{v_n}| \leq O(n^{2\kappa}\lambda_{\max}(\Sigma))) \geq 1 - p_n[\exp(-c_4n^{1-2\kappa}d_n^{-1}) + 6\exp(-c_5n^{1-2\kappa})];$$

if Condition H2 is satisfied and $d_n = o(n^{(1-2\kappa)/2})$, then

$$P(|\widehat{M}_{v_n}| \leq O(n^{2\kappa}\lambda_{\max}(\Sigma))) \geq 1 - p_n[\exp(-c_4n^{1-2\kappa}d_n^{-1}) + \exp(-c_5n^{1-2\kappa}d_n^{-2}) + 6\exp(-c_6n^{1-2\kappa})];$$

if Condition H3 is satisfied and $d_n = o(n^{(1-2\kappa)/3})$, then

$$P(|\widehat{M}_{v_n}| \leq O(n^{2\kappa}\lambda_{\max}(\Sigma))) \geq 1 - p_n[\exp(-c_4n^{1-2\kappa}d_n^{-1}) + 2\exp(-c_5n^{1/2-\kappa}d_n^{-3/2}) + 6\exp(-c_6n^{1-2\kappa})];$$

Type 3: if $d_n = o(n^{1-2\kappa})$, then

$$P(|\widehat{M}_{v_n}| \leq O(n^{2\kappa} \lambda_{\max}(\boldsymbol{\Sigma}))) \geq 1 - p_n [\exp(-c_4 n^{1-2\kappa} d_n^{-1}) + 12 \exp(-c_5 n^{1-2\kappa})].$$

The tail probabilities directly follow from Theorem 3.2 and have been explained in Section 3.3. The selected model size depends on the matrix $\boldsymbol{\Sigma}$ which involves the dependence structure of the covariates. As discussed in Fan, Feng and Song (2011), it can be assumed that $\lambda_{\max}(\boldsymbol{\Sigma}) = n^\tau$ for some $\tau > 0$. Correspondingly, the model size will be controlled at a reasonable rate of n . With iterative Goffins method described in the next Section 4, we will show in the simulation studies that the number of false positives can be very small while the true important variables are all selected even when the covariates are correlated.

3.5. *Connection and comparison with related literature.* When l is the squared error loss, Fan, Feng and Song's (2011) screening for nonparametric additive models is based on $n^{-1} \sum_{i=1}^n (\widehat{f}_{nj}^M(X_{i,j}))^2$. They presented a similar result to Theorem 3.2 here but with a different convergence probability as

$$P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n d_n \{ (8 + 2d_n) \exp(-c_4^* n^{1-4\kappa} d_n^{-3}) + 6d_n \exp(-c_5^* n d_n^{-3}) \},$$

and they can handle the NP dimensionality $\log p_n = o(n^{1-4\kappa} d_n^{-3})$. Compared with their result, our convergence probability in (20) for Gaussian regression not only improve on the convergence rate, but also improve significantly on those coefficient terms. We can achieve NP-dimensionality $\log p_n = o(n^{1-2\kappa} d_n^{-2})$. It should be noted that Fan, Feng and Song (2011) established the result under a weaker assumption than our Condition H2. More specifically, they assume $Y = \theta(\mathbf{X}) + \varepsilon$, $E(\varepsilon|\mathbf{X}) = 0$ and for any $B_1 > 0$, $E[\exp(B_1|\varepsilon)|\mathbf{X}] \leq B_2$ for some constant B_2 . On the other hand, this condition is stronger than our Condition H3. If we use (21) for the comparison here in favor of Fan, Feng and Song's (2011) result, then we can handle the NP-dimensionality $\log p_n = o(n^{1/2-\kappa} d_n^{-3/2})$. When $n^{1/3-2\kappa} = O(d_n)$, we can handle a higher dimensionality. Otherwise, their result is better.

When l is the quantile regression loss, He, Wang and Hong's (2013) screening is based on $n^{-1} \sum_{i=1}^n (\widehat{f}_{nj}^M(X_{i,j}))^2$. They have shown that for positive constants c_6^* , c_7^* and c_8^* ,

$$P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_6^* d_n n^{-4\kappa}) \leq 2 \exp(-c_7^* n^{1-8\kappa}) + \exp(-c_8^* n^{1-4\kappa} d_n^{-2}).$$

Correspondingly, they presented a convergence probability as

$$P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n \{ 11 \exp(-c_9^* n^{1-8\kappa}) + 12 d_n^2 \exp(-c_{10}^* n^{1-4\kappa} d_n^{-3}) \}.$$

Note that He, Wang and Hong (2013) consider a signal strength in Condition D as $c_1 n^{-2\kappa}$, and the parameter τ in Theorem 3.3 of He, Wang and Hong (2013) is equivalent to 2κ in our paper here. If we reset the minimum signal strength

in Condition D the same as that of He, Wang and Hong (2013), our result for Theorem 3.1 will be modified as

$$P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_3 d_n n^{-4\kappa}) \leq \exp(-c_4 n^{1-4\kappa} d_n^{-1})$$

under the conditions of He, Wang and Hong (2013). Correspondingly, our result for Theorem 3.2 will be modified as

$$P(M_\star \subset \widehat{M}_{v_n}) \geq 1 - s_n [\exp(-c_4 n^{1-4\kappa} d_n^{-1}) + 12 \exp(-c_5 n^{1-4\kappa})].$$

Therefore, He, Wang and Hong's (2013) convergence probability also indicates larger possibility of not selecting important covariates for a nonasymptotic setting.

When l is the negative log-likelihood loss for one-parameter exponential families, Fan and Song (2010) constructed sure screening for generalized linear models. If the B-spline approximation is treated as a type of group variable selection, then our Theorem 3.1 has some connections with Fan and Song's (2010) result. Compared with Fan and Song (2010), the tail probability in our Theorem 3.1 does not have the extra term $nP(\Omega_n^c)$ in their paper where n is the sample size and Ω_n is the region such that the loss function satisfies some Lipschitz condition. In Fan and Song (2010), their exponential bound also involves a Lipschitz constant. When the response variable is not bounded (e.g., most of the exponential families), this Lipschitz constant diverges to infinity, which results in a slower convergence rate for the tail bound, in contrast with our result. For example, when considering the squared error loss, Fan and Song (2010) Theorem 4 will have

$$P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\| \geq c_3 n^{-\kappa}) \leq \exp(-c_4 n^{(1-2\kappa)/3}) + nm_1 \exp(-c_4 n^{(1-2\kappa)/3})$$

for bounded covariates. For our Theorem 3.1 (under Condition H2), let $d_n = 2$ for a fair comparison, then we have

$$P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\| \geq c_3 n^{-\kappa}) \leq 2 \exp(-c_4 n^{1-2\kappa}).$$

It is clear that we have a much better result here. When considering the Poisson regression loss, the corresponding convergence rate in the tail probability bound can be much slower than $(1 - 2\kappa)/3$. For our Theorem 3.1 (under Condition H3), let $d_n = 2$ for a fair comparison, we have

$$P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\| \geq c_3 n^{-\kappa}) \leq \exp(-c_4 n^{1-2\kappa}) + 2 \exp(-c_5 n^{1/2-\kappa}).$$

It is still a better result than Fan and Song (2010).

4. Iterative Goffins method and improved variant. In practice, unimportant variables can be correlated with the important variables, therefore, such variables can have significant marginal effects even though they are not significant in the joint true model. To improve the performance of our screening method, we consider an iterative version of Goffins. Given the data $\{(\mathbf{X}_i, Y_i)\}, i = 1, \dots, n$, we choose the same truncation term $d_n = O(n^{1/5})$. In Theorem 3.2, the threshold v_n

is chosen at the level $d_n n^{-2\kappa}$. In practice, the parameter κ is unknown, but we can determine a data-driven threshold. To achieve this, we extend the random permutation idea of Fan, Feng and Song (2011) and Zhao and Li (2012). Let \mathbf{X} be the matrix with the i th row as \mathbf{X}_i . The algorithm works as follows:

Step 1. For every $j \in \{1, \dots, p\}$, compute

$$\hat{f}_{nj} = \arg \min_{f_{nj} \in S_n} \mathbb{P}_n l(f_{nj}(X_j), Y) \quad 1 \leq j \leq p.$$

Randomly permute the rows of \mathbf{X} and we have $\tilde{\mathbf{X}} = (\tilde{X}_1, \dots, \tilde{X}_p)$. Let $\omega_{(q)}$ be the q th quantile of $\{G_{n,j}^*, j = 1, \dots, p\}$, where $\hat{f}_{nj}^* = \arg \min_{f_{nj} \in S_n} \mathbb{P}_n l(f_{nj}(\tilde{X}_j), Y)$. Then our method selects the following variables: $\mathcal{A}_1 = \{j : G_{n,j} \geq \omega_{(q)}\}$. In our numerical studies, we choose $q = 1$, the maximum value of the empirical norm of the permuted estimates.

Step 2. Apply penalized regression on the set \mathcal{A}_1 to select a subset \mathcal{M}_1 . Specifically, when l is the negative log-likelihood loss, apply the penalized generalized additive model regression [e.g., penGAM in Meier, van de Geer and Bühlmann (2009)].

Step 3. For every $j \in \mathcal{M}_1^c = \{1, \dots, p\} / \mathcal{M}_1$, minimize $\mathbb{P}_n l(f_0 + \sum_{i \in \mathcal{M}_1} f_{ni}(X_i) + f_{nj}(X_j), Y)$ with respect to $f_0 \in \mathbb{R}$, $f_{ni} \in S_n$ for all $i \in \mathcal{M}_1$ and $f_{nj} \in S_n$. For identifiability, we apply the B spline basis without the intercept for $j \in \mathcal{M}_1^c$ and for $i \in \mathcal{M}_1$. Apply the screening procedure with adaptive threshold determined by the new random permutation. Choose a set of indices \mathcal{A}_2 . Then penalized regression is applied on the set $\mathcal{M}_1 \cup \mathcal{A}_2$ to select a subset \mathcal{M}_2 .

Step 4. Iterate the process until $|\mathcal{M}_l| \geq s_0$ or $\mathcal{M}_l = \mathcal{M}_{l-1}$.

This iterative version of Goffins will be denoted as ‘‘I-Goffins’’ in our simulation studies. To further stabilize the performance, we can apply a ‘‘cap’’ to control the number of selected variables in each iteration. For example, in our simulation studies, we restrict to select 1 variable at each step. Since the chance of selecting unimportant variables in each step has been reduced, the probability of selecting important variables in the subsequent steps has been improved. This is the idea behind the greedy INIS method proposed by Fan, Feng and Song (2011) for additive modeling. To be consistent with Fan, Feng and Song (2011), we name this improved variant of our method as greedy iterative goodness-of-fit nonparametric screening (GI-Goffins).

5. Simulation studies. Similar to Fan, Feng and Song (2011), we set $n = 400$ but we consider $p = 1000, 2000, 5000$ for all examples to investigate the impact of high dimensionality on screening methods. Following Fan, Feng and Song (2011), we consider the number of spline basis functions as $d_n = \lceil n^{1/5} \rceil + 2 = 6$. Note that in this paper we consider the full B spline basis, and Fan, Feng and Song (2011) considered the B spline basis without the intercept. The goodness-of-fit screening methods under the two sets of basis are equivalent. Eight simulation examples

will be constructed according to the four major types of regressions: Gaussian regression, Logistic regression, Poisson regression and quantile regression. Define

$$f_1(x) = x, \quad f_2(x) = (2x - 1)^2, \quad f_3(x) = \sin(2\pi x)/(2 - \sin(2\pi x)),$$

$$f_4(x) = 0.1 \sin(2\pi x) + 0.2 \cos(2\pi x) + 0.3 \sin(2\pi x)^2 \\ + 0.4 \cos(2\pi x)^3 + 0.5 \sin(2\pi x)^3,$$

$$f_5(x) = \exp(x - 0.5), \quad f_6(x) = 0.1 \sin(2\pi x)^2 + 0.4 \cos(2\pi x)^3,$$

$$f_7(x) = \sin(x - 1), \quad f_8(x) = (x - 1.5)^2, \quad f_9(x) = 2 \cos(x)/(2 - \sin(x)).$$

- Model 1 (Linear regression): $Y|\mathbf{X} = 5f_1(X_1) + 3f_2(X_2) + 4f_3(X_3) + 6f_4(X_4) + \sqrt{1.74}\varepsilon$. Each $X_i \sim \text{Uniform}(0, 1)$ i.i.d. and $\varepsilon \sim N(0, 1)$.
- Model 2 (Linear regression): The model is the same as Model 1 but the covariates $\mathbf{X} = (X_1, \dots, X_p)^T$ are simulated according to the random effects model $X_j = (W_j + tU)/(1 + t)$, $j = 1, \dots, p$ where W_1, \dots, W_p and U are i.i.d. $\text{Uniform}(0, 1)$ and $t = 0.4$.
- Model 3 (Logistic regression): $\ln(P(Y = 1|\mathbf{X})/P(Y = 0|\mathbf{X})) = 2f_1(X_1) + 3f_7(X_2) + 2f_8(X_3) + 3.5f_9(X_4)$. Each $X_i \sim \text{Uniform}(-2.5, 2.5)$ i.i.d.
- Model 4 (Logistic regression): The model is the same as Model 3 but the covariates $\mathbf{X} = (X_1, \dots, X_p)^T$ are simulated according to the random effects model $X_j = (W_j + tU)/(1 + t)$, $j = 1, \dots, p$ where W_1, \dots, W_p are i.i.d. from $\text{Uniform}(-2.5, 2.5)$, independent of $U \sim \text{Uniform}(0, 1)$ and $t = 0.4$.
- Model 5 (Poisson regression): $Y|\mathbf{X} \sim \text{Poisson}(\exp\{f_1(X_1) + f_3(X_2) + f_5(X_3) + f_6(X_4)\})$. Each $X_i \sim \text{Uniform}(0, 1)$ i.i.d.
- Model 6 (Poisson regression): The model is the same as Model 5 and the covariates $\mathbf{X} = (X_1, \dots, X_p)^T$ are simulated according to the same structure as Model 2.
- Model 7 (Heteroscedastic regression): $Y|\mathbf{X} = 5f_1(X_1) + 3f_2(X_2) + 4f_3(X_3) + 4f_5(X_4) + 0.5 \exp(f_6(X_{20}) + f_7(X_{21}) + f_8(X_{22}))\varepsilon$, where $\mathbf{X} \sim N_p(0, \Sigma)$ independent of $\varepsilon \sim N(0, 1)$ and the (i, j) th element of covariance matrix Σ is $0.8^{|i-j|}$.
- Model 8 (Heteroscedastic regression): The model is the same as Model 7 except that the random error $\varepsilon \sim \text{Laplace}(0, 2)$.

Models 1 and 2 have been similarly considered in Meier, van de Geer and Bühlmann (2009) and Fan, Feng and Song (2011), while Models 3–8 are newly proposed in the current paper. The covariates are independent in Models 1, 3, 5 but correlated in Models 2, 4, 6, 7, 8. Note that Model 8 is different from Model 7 because the $\text{Laplace}(0, 2)$ distribution for random error will emphasize more on the covariates X_{20}, X_{21}, X_{22} , thus making the heteroscedastic regression model more challenging.

Minimum model size. Following Fan and Song (2010), Fan, Feng and Song (2011) as well as later literature in sure screening field, we use the minimum model

size required to contain the true model M_* as a measure of the effectiveness of a screening method. The simulation round is 500 for all the examples. We compare our Goffins method with six other successful screening methods in the existing literature, including some recent model-free screening methods. More specifically, we consider fused Kolmogorov filter (Kfilter) by Mai and Zou (2015), quantile adaptive screening (QaSIS) by He, Wang and Hong (2013), SIS for generalized linear model by Fan and Song (2010), sure independent ranking and screening (SIRS) by Zhu, Li, Li and Zhu (2012), distance correlation learning (DC) by Li, Zhong and Zhu (2012) and empirical likelihood screening (EL) by Chang, Tang and Wu (2013). Note that when considering squared error loss, our Goffins is equivalent to NIS by Fan, Feng and Song (2011). Therefore, we will treat NIS as a special example of Goffins, and will not present NIS as a separate method for comparison here. However, when considering quantile regression loss, our Goffins method is different from QaSIS in He, Wang and Hong (2013), because Goffins is based on goodness-of-fit statistics while QaSIS is based on squared norm of fitted nonparametric function. We choose quantile 75% whenever a quantile regression loss is considered. The implementations of QaSIS and SIRS are based on <http://users.stat.umn.edu/~wangx346/research/example1b.txt>. The implementations of Kfilter, DC and EL are based on the R codes from the authors of related literature. The implementation of SIS is based on the R package “SIS.”

In Tables 1–2 along with the Tables S1–S2 in the Supplementary Material [Han (2019)], we present the median, the interquartile range (IQR) and different quantiles of minimum model size. Following existing literature, if the median is closer to the true model size and the IQR is smaller, the corresponding screening method is considered as more effective. Overall, our Goffins method performs best among the seven screening methods. For Models 1, 3, 4, 5, 6, our medians of minimum model size are close to the true model size 4 and IQRs are the smallest. Our medians and IQRs will not increase significantly when the dimensionality p increases from 1000 to 5000. For comparison, other methods tend to select a much larger model to contain the true model, and the performance can deteriorate dramatically when p increases. Furthermore, 5%, 25%, 75% and 95% quantiles of our minimum model size are significantly smaller than the other methods. Models 7–8 are very challenging heteroscedastic regression models, but our method still performs better than the other methods, including Kfilter and QaSIS. Table S2 in the Supplementary Material [Han (2019)] also suggests that even when we consider quantile regression loss, Goffins is different from QaSIS. Model 2 turns out to be a difficult example for all the methods. However, our simulation in tables S3–S5 of Supplementary Material [Han (2019)] will show that an iterative version of Goffins (GI-Goffins or I-Goffins) can substantially reduce the false positives while selecting the true important variables.

TABLE 1
The median (IQR) and 5%, 25%, 75%, 95% quantiles of minimum model size

| Model | p | Methods | Median (IQR) | 5% | 25% | 75% | 95% | |
|---------------------|----------------|---------------------|------------------|-----------|---------|---------|---------|--------|
| Model 3 Size = 4 | 1000 | Goffins | 5 (7.25) | 4 | 4 | 11.25 | 67.05 | |
| | | Kfilter | 14 (28) | 4 | 6 | 34 | 115.05 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 393.5 (537) | 38.9 | 158 | 695 | 941.05 | |
| | | SIRS | 393 (537) | 38.95 | 158 | 695 | 941.05 | |
| | | DC | 16 (25) | 4 | 9 | 34 | 112.05 | |
| | | EL | 451 (524.25) | 55.75 | 211 | 725.25 | 948.05 | |
| | | 2000 | Goffins | 5 (6) | 4 | 4 | 10 | 64.05 |
| | Kfilter | | 13 (28.25) | 4 | 7 | 35.25 | 119.65 | |
| | QaSIS | | NA | NA | NA | NA | NA | |
| | SIS | | 376.5 (521) | 29.85 | 168.50 | 689.5 | 954.05 | |
| | SIRS | | 376 (521.5) | 29.85 | 168 | 689.5 | 954.05 | |
| | DC | | 16 (24.25) | 4 | 8 | 32.25 | 105.20 | |
| | EL | | 441 (508.25) | 38.95 | 215.50 | 723.75 | 961.05 | |
| | 5000 | | Goffins | 9 (36.25) | 4 | 4 | 40.25 | 301.20 |
| | | Kfilter | 59 (151.5) | 5 | 20 | 171.5 | 597 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 2164 (2819.5) | 165.1 | 887.5 | 3707.0 | 4814.4 | |
| | | SIRS | 2163.5 (2817.5) | 164.20 | 888.75 | 3706.25 | 4814.40 | |
| | | DC | 67 (130.5) | 8 | 27.75 | 158.25 | 504.15 | |
| | | EL | 2465.5 (2705.75) | 229.85 | 1136.25 | 3842.00 | 4846.05 | |
| | | Model 4 Size = 4 | 1000 | Goffins | 5 (6) | 4 | 4 | 10 |
| | Kfilter | | | 16 (36) | 4 | 6 | 42 | 163.1 |
| | QaSIS | | | NA | NA | NA | NA | NA |
| SIS | 418.5 (566) | | | 20.9 | 148.75 | 714.75 | 920.3 | |
| SIRS | 418.5 (566.25) | | | 19.95 | 148.50 | 714.75 | 920.3 | |
| DC | 17 (33) | | | 4 | 8 | 41 | 137.05 | |
| EL | 495 (543.75) | | | 42 | 211.75 | 755.50 | 933.15 | |
| 2000 | Goffins | | | 6 (17) | 4 | 4 | 21 | 131.1 |
| | Kfilter | | 32.5 (72.25) | 4 | 10 | 82.25 | 243.35 | |
| | QaSIS | | NA | NA | NA | NA | NA | |
| | SIS | | 858 (1079.25) | 36.85 | 312.25 | 1391.50 | 1892 | |
| | SIRS | | 858 (1077) | 36.9 | 314.5 | 1391.5 | 1892 | |
| | DC | | 34 (70) | 5 | 14 | 84 | 240.15 | |
| | EL | | 1008 (1037.25) | 83.85 | 446.75 | 1484 | 1911.05 | |
| | 5000 | | Goffins | 9 (32) | 4 | 4 | 36 | 332.4 |
| Kfilter | | | 59.5 (174) | 5 | 16 | 190 | 770.75 | |
| QaSIS | | | NA | NA | NA | NA | NA | |
| SIS | | | 1857.5 (2778.5) | 84.95 | 749.50 | 3528.00 | 4675.10 | |
| SIRS | | | 1855.5 (2781.5) | 84.95 | 746.50 | 3528.00 | 4675.10 | |
| DC | | | 70 (160.5) | 8 | 26 | 186.5 | 661.0 | |
| EL | | | 2261.5 (2701) | 196.8 | 1054.5 | 3755.5 | 4726.5 | |

TABLE 2
The median (IQR) and 5%, 25%, 75%, 95% quantiles of minimum model size

| Model | p | Methods | Median (IQR) | 5% | 25% | 75% | 95% | |
|---------------------|-----------------------|----------------|---------------------|-----------|------------|------------|------------|--------|
| Model 5 Size = 4 | 1000 | Goffins | 4 (0) | 4 | 4 | 4 | 12.1 | |
| | | Kfilter | 28 (63) | 4 | 11 | 74 | 265.05 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 447 (525.75) | 22.95 | 174.25 | 700 | 944.15 | |
| | | SIRS | 491 (501) | 46 | 242.5 | 743.5 | 951.05 | |
| | | DC | 11 (18) | 4 | 6 | 24 | 71.05 | |
| | | EL | 461.5 (510) | 32.9 | 194.75 | 704.75 | 949 | |
| | 2000 | Goffins | 4 (1) | 4 | 4 | 5 | 17 | |
| | | Kfilter | 60 (144) | 5 | 17 | 161 | 504.1 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 980 (1154.25) | 48.95 | 371.75 | 1526 | 1924.05 | |
| | | SIRS | 1031 (1028) | 112.75 | 516 | 1544 | 1904.30 | |
| | | DC | 21 (36) | 4 | 9 | 45 | 132 | |
| | | EL | 1013 (1166.75) | 54.9 | 393.75 | 1560.5 | 1923.05 | |
| | 5000 | Goffins | 4 (2) | 4 | 4 | 6 | 42 | |
| | | Kfilter | 147.5 (377) | 7 | 41.75 | 418.75 | 1340.00 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 2237.5 (2662.75) | 163.80 | 978.25 | 3641.00 | 4748.45 | |
| | | SIRS | 2619.5 (2518.75) | 266.80 | 1468.75 | 3987.50 | 4723.20 | |
| | | DC | 44 (98.25) | 6 | 17 | 115.25 | 360.70 | |
| | | EL | 2320 (2701) | 183.85 | 988.00 | 3689.00 | 4773.15 | |
| | Model 6 Size = 4 | 1000 | Goffins | 4 (1.25) | 4 | 4 | 5.25 | 19.05 |
| | | | Kfilter | 35.5 (95) | 4.95 | 12 | 107 | 345.05 |
| | | | QaSIS | NA | NA | NA | NA | NA |
| SIS | | | 400.5 (571.75) | 28 | 150.25 | 722 | 940.05 | |
| SIRS | | | 510.5 (514) | 49.95 | 244.5 | 758.5 | 952.15 | |
| DC | | | 17 (36) | 4 | 7 | 43 | 129 | |
| EL | | | 425 (567) | 30.95 | 163 | 730 | 940.05 | |
| 2000 | | Goffins | 4 (2) | 4 | 4 | 6 | 31 | |
| | | Kfilter | 66.5 (159) | 5 | 22 | 181 | 552.7 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 930.5 (1067.25) | 41.95 | 356 | 1423.25 | 1905.10 | |
| | | SIRS | 1027.5 (1051.75) | 94.80 | 488.50 | 1540.25 | 1911.05 | |
| | | DC | 23 (47) | 5 | 10 | 57 | 240.05 | |
| | | EL | 972.5 (1059.5) | 50.80 | 383.25 | 1442.75 | 1912 | |
| 5000 | | Goffins | 5 (7) | 4 | 4 | 11 | 74.4 | |
| | | Kfilter | 182 (483.25) | 6 | 52 | 535.25 | 1614.10 | |
| | | QaSIS | NA | NA | NA | NA | NA | |
| | | SIS | 2072.5 (2690.75) | 198 | 825 | 3515.75 | 4709.70 | |
| | | SIRS | 2442 (2528.25) | 231.70 | 1324 | 3852.25 | 4778.45 | |
| | | DC | 74 (173.5) | 6 | 23.75 | 197.25 | 593.15 | |
| | | EL | 2130 (2627.5) | 214.6 | 926.5 | 3554.0 | 4710.3 | |

6. Data analysis. Classification between the malignant pleural mesothelioma (MPM) and the lung cancer adenocarcinoma (ADCA) has received increasing attention in both clinical studies and high dimensional statistical research. Gordon et al. (2002) studied the data from 181 tissue samples (31 MPM and 151 ADCA) with 12,533 gene expression levels for each sample. Among these 181 sample data, 16 MPM and 16 ADCA have been combined as the training set while the other 149 samples (15 MPM and 134 ADCA) are considered as the testing set. The goal of research is in two-fold as explained in Gordon et al. (2002): 1. Find the minimum number of predictor genes that are most importantly associated with the disease type; 2. Construct a classifier rule which can predict the future patients' disease type based on their gene expression levels with high statistical accuracy. Aspect 1 can substantially reduce the medical cost of obtaining patients' relevant gene data and the cost of potential scientific experiments on such genes. The performance of the classifier is usually evaluated based on the testing data.

Since the disease type is a categorical data, and the number of genes is extremely high ($p = 12,533$) compared with the small sample size ($n = 32$), we will apply our GI-Goffins method with respect to the logistic regression for the training data. We first standardize the gene expression data for each gene over the training samples such that the sample mean is 0 and the sample standard deviation is 1. Our method selects five genes that are importantly associated with the disease type: "31575-f-at," "37716-at," "39795-at," "41286-at" and "41402-at." We construct a generalized additive model (B spline basis without the intercept and the number of spline basis functions as $d_n = \lceil n^{1/5} \rceil + 1 = 3$) based on such five genes and apply the model to the training data. The fitted nonparametric functions corresponding to those five genes have also been plotted in Figure S1 in the Supplementary Material [Han (2019)]. Then we apply our constructed model to the test data. Among the 149 samples for the testing data, we make 144 correct predictions. For the 5 samples that we misclassified, one MPM sample has been predicted as ADCA while four ADCA samples has been predicted as MPM. ISIS for the generalized linear model has also been considered to select important variables. To be fair, we also apply a generalized additive model based on the selected genes for the training data and further use this fitted model for classification on the test data. However, this method will select fewer and different genes and the performance is much inferior to our method. I-EL is an iterative version of EL and penalized empirical likelihood regression described in Chang, Tang and Wu (2013). Its performance is even worse than ISIS.

This lung cancer data has also been analyzed by various statistical methods in the past literature. It is impossible and unnecessary for us to list all the relevant results here, and we only compare our method with some representative methods which have been shown superior performance. In Table 3, we will compare our GI-Goffins method with linear discriminant methods such as ROAD in Fan, Feng and Tong (2012) and FAIR in Fan and Fan (2008). Our GI-Goffins is a good balance between the testing error and the number of selected genes compared with other

TABLE 3
Performance of methods on lung cancer data. $p = 12,533$

| Method | Training error | Testing error | Number of selected genes |
|------------|----------------|---------------|--------------------------|
| GI-Goffins | 0/32 | 5/149 | 5 |
| ISIS | 0/32 | 18/149 | 2 |
| I-EL | 0/32 | 40/149 | 5 |
| ROAD | 1/32 | 1/149 | 52 |
| FAIR | 0/32 | 7/149 | 31 |

methods. More selected genes will cause substantial cost in future diagnosis and experiments. Therefore, GI-Goffins is the method that we recommend for practice.

7. Further discussions.

7.1. *Optimality.* An interesting question is whether the convergence rate in the upper bound of the tail probability that we established in Theorem 3.1 is optimal. More specifically, if we have

$$b_1 \exp(-c_2 n^a) \leq P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_1 d_n n^{-2\kappa}) \leq b_2 \exp(-c_3 n^a)$$

for some constants a, b_1, b_2, c_1, c_2 and c_3 , then we can say that the convergence rate a in the upper bound of the tail probability is optimal, because the convergence rate a cannot be improved further.

When the loss function l is the negative log likelihood loss of one-parameter exponential families, under general regularity conditions, the maximum likelihood estimator has the asymptotic normality [Heyde (1997), Gao et al. (2008)], that is,

$$[I_j(\boldsymbol{\beta}_j^M)]^{1/2}(\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M) - N(0, \mathbf{I}_{d_n}) \rightarrow 0 \quad \text{in distribution,}$$

where I_j is the information matrix of the j th covariate. Plugging in the negative log likelihood loss and the B-spline basis functions, we have

$$n^{1/2}\{E[b''(\boldsymbol{\Psi}_j^T \boldsymbol{\beta}_j^M) \boldsymbol{\Psi}_j \boldsymbol{\Psi}_j^T]\}^{1/2}(\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M) - N(0, \mathbf{I}_{d_n}) \rightarrow 0 \quad \text{in distribution.}$$

Since $\boldsymbol{\Psi}_j^T \boldsymbol{\beta}_j^M$ is bounded based on the argument in the Supplementary Material [Han (2019)] and $b''(\cdot)$ is a continuous function, $b''(\boldsymbol{\Psi}_j^T \boldsymbol{\beta}_j^M)$ is upper bounded by a positive constant. Furthermore, due to Lemma 3 in the Supplementary Material [Han (2019)],

$$(\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M)^T n \{E[b''(\boldsymbol{\Psi}_j^T \boldsymbol{\beta}_j^M) \boldsymbol{\Psi}_j \boldsymbol{\Psi}_j^T]\}(\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M) \leq D_2 n d_n^{-1} \|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2.$$

Therefore, asymptotically, we have

$$\begin{aligned} P(\|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq c_1 d_n n^{-2\kappa}) &= P(D_2 n d_n^{-1} \|\widehat{\boldsymbol{\beta}}_j^M - \boldsymbol{\beta}_j^M\|^2 \geq D_2 c_1 n^{1-2\kappa}) \\ &= P(\chi_{d_n}^2 \geq D_2 c_1 n^{1-2\kappa}). \end{aligned}$$

Thus, we need to find a lower bound for the tail probability of $\chi_{d_n}^2$ distribution. When $d_n = 1$, it is well known that for any positive y ,

$$P(\chi_1^2 \geq y) \geq 1 - \sqrt{1 - \exp\left(-\frac{2y}{\pi}\right)} = \frac{\exp\left(-\frac{2y}{\pi}\right)}{1 + \sqrt{1 - \exp\left(-\frac{2y}{\pi}\right)}} \geq \frac{1}{2} \exp\left(-\frac{2y}{\pi}\right).$$

Let $y = D_2 c_1 n^{1-2\kappa}$, comparing with our Theorem 3.1 under Condition H1 or H2, we have achieved the optimal convergence rate $n^{1-2\kappa}$ asymptotically. When $d_n = 2$, for any positive y , $P(\chi_2^2 \geq y) = \exp(-\frac{y}{2})$. Comparing with our Theorem 3.1 under Condition H1 or H2, we have also achieved the optimal convergence rate $n^{1-2\kappa}$ asymptotically. For more general d_n , we do not have a sharp lower bound of the tail probability of Chi-square distribution. Therefore, we will not discuss further here.

7.2. Adaptive threshold. Theorem 3.2 is established based on a threshold v_n at the level of $d_n n^{-2\kappa}$. In practice, the parameter κ is unknown. Therefore, we need an adaptive threshold for the real data. Consider a threshold \widehat{v}_n which is constructed based on the sample data, it will be interesting to derive a lower bound for $P(M_\star \subseteq \widehat{M}_{\widehat{v}_n})$. We have

$$P(M_\star \subseteq \widehat{M}_{\widehat{v}_n}) \geq 1 - \sum_{j \in M_\star} P(G_{n,j} < \widehat{v}_n).$$

Note that

$$P(G_{n,j} < \widehat{v}_n) \leq P(G_{n,j} < v_n) + P(\widehat{v}_n > v_n).$$

We have derived the upper bound of $P(G_{n,j} < v_n)$ is the proof of Theorem 3.2. Therefore, we need to derive an upper bound for the second term here.

Consider a permutation of the sample covariates $\{\mathbf{X}_i\}_{i=1}^n$. We can obtain the estimates of marginal regression based on the permuted data:

$$(\widehat{\boldsymbol{\beta}}_j^M)^\pi = \arg \min_{\boldsymbol{\beta}_j \in \mathbb{R}^{d_n}} \mathbb{P}_n l(\boldsymbol{\Psi}_j^T(X_j^\pi) \boldsymbol{\beta}_j, Y),$$

where $\pi = (\pi_1, \dots, \pi_n)$ is a permutation of the index $\{1, 2, \dots, n\}$. Note that $(G_{n,j})^\pi$ is a statistical estimate of 0. We will derive an upper bound for $P(\widehat{v}_n > v)$ for a special case where the loss function $l(\cdot)$ is the squared error loss. Note that the least squares estimate follows

$$\sqrt{n} \widehat{\boldsymbol{\beta}}_j^\pi = \left(\frac{1}{n} \boldsymbol{\Psi}(X_j^\pi) \boldsymbol{\Psi}(X_j^\pi)^T \right)^{-1} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \boldsymbol{\Psi}^T(X_j^{\pi_i}) Y_i \right).$$

By Anderson and Robinson (2001) Theorem 3.2, after some algebra, we have the asymptotic normality

$$\sqrt{n} \widehat{\boldsymbol{\beta}}_j^\pi - N(0, \widetilde{\boldsymbol{\Sigma}}) \rightarrow 0 \quad \text{in distribution,}$$

where $\widetilde{\boldsymbol{\Sigma}} = (E \boldsymbol{\Psi} \boldsymbol{\Psi}^T)^{-1} [E(\boldsymbol{\Psi} - E \boldsymbol{\Psi})(\boldsymbol{\Psi} - E \boldsymbol{\Psi})^T] (E \boldsymbol{\Psi} \boldsymbol{\Psi}^T)^{-1}$.

If we let $\widehat{v}_{n,j} = \frac{1}{n} \sum_{i=1}^n (\Psi^T(X_j^{\pi_i}) \widehat{\beta}_j^\pi)^2$, different thresholds for the marginal utilities of different covariates. As we have discussed in the Section 2.2, this screening based on $\mathbb{P}_n(\Psi_j^T \widehat{\beta}_j^M)^2$ is equivalent to our goodness-of-fit screening when the loss function is the squared error loss. We can show that asymptotically

$$P(\widehat{v}_{n,j} \geq c_1 d_n n^{-2\kappa}) \leq P(\chi_{d_n}^2 \geq c_2 d_n n^{-2\kappa}) \leq \exp(-c_3 d_n n^{1-2\kappa})$$

for some constants c_1, c_2 and c_3 . The second inequality is by [Laurent and Massart \(2000\)](#) Lemma 3.1. Correspondingly, we have a lower bound for the convergence probability of containing the true model. In practice, the threshold \widehat{v}_n will be chosen as the maximum value of $\{\widehat{v}_{n,j}\}_{j=1}^{p_n}$ under a number of permutations, and the loss function can be general. We do not have a theoretical result for such more complicated situations.

7.3. Choice of loss function. The framework of goodness-of-fit nonparametric screening includes many screening methods based on the choice of loss functions. An important question is how to choose loss function for practical data. When the response variable Y takes values $\{0, 1\}$, we suggest to consider the logistic regression loss: $l(\omega, Y) = -\omega Y + \ln(1 + \exp(\omega))$; when Y takes nonnegative integer values, we suggest to consider the Poisson regression loss: $l(\omega, Y) = -Y\omega + \exp(\omega) + \ln(Y!)$; When the distribution of Y is expected to be complicated, the quantile regression loss can be considered; when Y is a continuous variable, we suggest to start with the Gaussian regression loss: $l(\omega, Y) = (Y - \omega)^2/2$. This brief guideline could raise misspecification issue of loss functions.

7.4. Iterative screening procedure. The idea of iterative screening and penalization has been proposed since [Fan and Lv \(2008\)](#), and has achieved numerical success in practice. However, formal theoretical justification is still an open problem in the field. The first step is a marginal screening. To simplify the discussion, assume a fixed threshold γ_n is applied the selected variables $A_1 = \{j : G_{n,j} \geq \gamma_n\}$ satisfies $M_\star \subseteq A_1$ with high probability (sure screening property). For the second step, based on the set A_1 , we apply some penalized regression and select a subset M_1 . Ideally, we want to show sign consistency for M_1 under some regularity conditions. The difficulty is that the set A_1 is random, which is different from a conventional penalization regression. Fortunately, we can borrow the technique in [Weng, Feng and Qiao \(2017\)](#), which considers a two-step procedure for linear regression model (similar to screening + penalization). For the third step, it is a conditional marginal screening after penalization. [Barut, Fan and Verhasselt \(2016\)](#) has shown the sure screening property based on the conditional screening for generalized linear model. Therefore, if the sign consistency is achieved in step 2, then under some regularity conditions, sure screening property can be achieved in step 3. By mathematical induction, the iterative procedure can achieve sign consistency. We would like to explore the technical details as our future studies.

Acknowledgments. The author wants to thank the Joint Editor, the Associate Editor and the three anonymous referees for many insightful comments which significantly improve the presentation of the paper.

The author deeply appreciates Professor Jianqing Fan for his encouragement and constructive comments on this project. The author would like to thank Dr. Cheng Yong Tang and Professor Linda Zhao for helpful discussions on the paper. The author also thanks Dr. Shujie Ma and Dr. Lily Wang for the helpful discussion on B spline basis.

Special thanks go to the following researchers who kindly share their codes of numerical studies: Dr. Yang Feng for NIS, Dr. Qing Mai for Kfilter, Dr. Lukas Meier for penGAM, Dr. Cheng Yong Tang for EL and Dr. Wei Zhong for DC.

SUPPLEMENTARY MATERIAL

Supplement to “Nonparametric screening under conditional strictly convex loss for ultrahigh dimensional sparse data” (DOI: [10.1214/18-AOS1738SUPP](https://doi.org/10.1214/18-AOS1738SUPP); .pdf). Due to the space limit, all the technical proofs as well as some numerical results are relegated to the Supplementary Material [Han (2019)].

REFERENCES

- ANDERSON, M. J. and ROBINSON, J. (2001). Permutation tests for linear models. *Aust. N. Z. J. Stat.* **43** 75–88. [MR1837497](#)
- BARUT, E., FAN, J. and VERHASSELT, A. (2016). Conditional sure independence screening. *J. Amer. Statist. Assoc.* **111** 1266–1277. [MR3561948](#)
- BRÈGMAN, L. M. (1967). A relaxation method of finding a common point of convex sets and its application to the solution of problems in convex programming. *Ž. Vychisl. Mat. Mat. Fiz.* **7** 620–631. [MR0215617](#)
- mr BULDYGIN, V. and KOZACHENKO, Y. (2000). Metric characterization of random variables and random processes.. *Translations of Mathematical Monographs* **188**.
- CANDÈS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *Ann. Statist.* **35** 2313–2404.
- CHANG, J., TANG, C. Y. and WU, Y. (2013). Marginal empirical likelihood and sure independence feature screening. *Ann. Statist.* **41** 2123–2148. [MR3127860](#)
- DE BOOR, C. (1978). *A Practical Guide to Splines. Applied Mathematical Sciences* **27**. Springer, New York. [MR0507062](#)
- FAN, J. and FAN, Y. (2008). High-dimensional classification using features annealed independence rules. *Ann. Statist.* **36** 2605–2637. [MR2485009](#)
- FAN, J., FENG, Y. and SONG, R. (2011). Nonparametric independence screening in sparse ultrahigh-dimensional additive models. *J. Amer. Statist. Assoc.* **106** 544–557. [MR2847969](#)
- FAN, J., FENG, Y. and TONG, X. (2012). A road to classification in high dimensional space: The regularized optimal affine discriminant. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **74** 745–771. [MR2965958](#)
- FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.* **96** 1348–1360. [MR1946581](#)
- FAN, J. and LV, J. (2008). Sure independence screening for ultrahigh dimensional feature space. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **70** 849–911. [MR2530322](#)

- FAN, J., MA, Y. and DAI, W. (2014). Nonparametric independence screening in sparse ultra-high-dimensional varying coefficient models. *J. Amer. Statist. Assoc.* **109** 1270–1284. [MR3265696](#)
- FAN, J., SAMWORTH, R. and WU, Y. (2009). Ultrahigh dimensional feature selection: Beyond the linear model. *J. Mach. Learn. Res.* **10** 2013–2038. [MR2550099](#)
- FAN, J. and SONG, R. (2010). Sure independence screening in generalized linear models with NP-dimensionality. *Ann. Statist.* **38** 3567–3604. [MR2766861](#)
- FREUND, Y. and SCHAPIRE, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. System Sci.* **55** 119–139. [MR1473055](#)
- GAO, Q., WU, Y., ZHU, C. and WANG, Z. (2008). Asymptotic normality of maximum quasi-likelihood estimators in generalized linear models with fixed design. *J. Syst. Sci. Complex.* **21** 463–473. [MR2425677](#)
- GORDON, G. et al. (2002). Translation of Microarray Data into Clinically Relevant Cancer Diagnostic Tests Using Gene Expression Ratios in Lung Cancer and Mesothelioma. *Cancer Research* **62** 4963–4967.
- HAN, X. (2019). Supplement to “Nonparametric screening under conditional strictly convex loss for ultrahigh dimensional sparse data.” DOI:10.1214/18-AOS1738SUPP.
- HE, X., WANG, L. and HONG, H. G. (2013). Quantile-adaptive model-free variable screening for high-dimensional heterogeneous data. *Ann. Statist.* **41** 342–369. [MR3059421](#)
- HEYDE, C. C. (1997). *Quasi-Likelihood and Its Application: A General Approach to Optimal Parameter Estimation*. Springer, New York. [MR1461808](#)
- KOENKER, R. (2005). *Quantile Regression. Econometric Society Monographs* **38**. Cambridge Univ. Press, Cambridge. [MR2268657](#)
- LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28** 1302–1338. [MR1805785](#)
- LI, R., ZHONG, W. and ZHU, L. (2012). Feature screening via distance correlation learning. *J. Amer. Statist. Assoc.* **107** 1129–1139. [MR3010900](#)
- LI, G., PENG, H., ZHANG, J. and ZHU, L. (2012). Robust rank correlation based screening. *Ann. Statist.* **40** 1846–1877. [MR3015046](#)
- MAI, Q. and ZOU, H. (2015). The fused Kolmogorov filter: A nonparametric model-free screening method. *Ann. Statist.* **43** 1471–1497. [MR3357868](#)
- MEIER, L., VAN DE GEER, S. and BÜHLMANN, P. (2009). High-dimensional additive modeling. *Ann. Statist.* **37** 3779–3821. [MR2572443](#)
- SONG, R., LU, W., MA, S. and JENG, X. J. (2014). Censored rank independence screening for high-dimensional survival data. *Biometrika* **101** 799–814. [MR3286918](#)
- STONE, C. J. (1986). The dimensionality reduction principle for generalized additive models. *Ann. Statist.* **14** 590–606. [MR0840516](#)
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. [MR1379242](#)
- WENG, H., FENG, Y. and QIAO, X. (2017). Regularization after retention in ultrahigh dimensional linear regression models. *Statist. Sinica*. In press.
- ZHANG, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.* **38** 894–942. [MR2604701](#)
- ZHANG, C., JIANG, Y. and SHANG, Z. (2009). New aspects of Bregman divergence in regression and classification with parametric and nonparametric estimation. *Canad. J. Statist.* **37** 119–139. [MR2509465](#)
- ZHAO, S. D. and LI, Y. (2012). Principled sure independence screening for Cox models with ultrahigh-dimensional covariates. *J. Multivariate Anal.* **105** 397–411. [MR2877525](#)
- ZHU, L.-P., LI, L., LI, R. and ZHU, L.-X. (2011). Model-free feature screening for ultrahigh-dimensional data. *J. Amer. Statist. Assoc.* **106** 1464–1475. [MR2896849](#)
- ZOU, H. (2006). The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* **101** 1418–1429. [MR2279469](#)

- ZOU, H. and HASTIE, T. (2005). Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **67** 301–320. [MR2137327](#)
- ZOU, H. and LI, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models. *Ann. Statist.* **36** 1509–1533. [MR2435443](#)

DEPARTMENT OF STATISTICAL SCIENCE
FOX BUSINESS SCHOOL
TEMPLE UNIVERSITY
PHILADELPHIA, PENNSYLVANIA 19122
USA
E-MAIL: hanxu3@temple.edu