Statistica Sinica Preprint No: SS-2020-0303							
Title	Model Checking in Large-Scale Dataset via						
	Structure-Adaptive-Sampling						
Manuscript ID	SS-2020-0303						
URL	http://www.stat.sinica.edu.tw/statistica/						
DOI	10.5705/ss.202020.0303						
Complete List of Authors	Yixin Han,						
	Ping Ma,						
	Haojie Ren and						
	Zhaojun Wang						
Corresponding Author	Zhaojun Wang						
E-mail	zjwang@nankai.edu.cn						
Notice: Accepted version subje	ct to English editing.						

MODEL CHECKING IN LARGE-SCALE DATASET VIA STRUCTURE-ADAPTIVE-SAMPLING

Yixin Han¹, Ping Ma², Haojie Ren³ and Zhaojun Wang¹

¹School of Statistics and Data Science, LPMC & KLMDASR, Nankai University, Tianjin, P.R. China

²Department of Statistics, University of Georgia, Athens, GA, USA

³School of Mathematical Sciences, Shanghai Jiao Tong University, Shanghai, P.R. China

Abstract: Lack-of-fit testing is often essential in many applications of statistical/machine learning. Despite the availability of large-scale datasets, the challenges associated with model checking when some resource budgets are limited are not yet well addressed. In this paper, we propose a design-adaptive testing procedure to check a general model when only a limited number of data observations are available. We derive an optimal sampling strategy to select a small subset from a large pool of data, Structure-Adaptive-Sampling, with which the proposed test possesses the asymptotically best power. Numerical results on both synthetic and real-world data confirm the effectiveness of the proposed method.

Keywords and phrases: Dimension reduction; Kernel smoothing; Large-scale dataset; Nonparametric lack-of-fit tests; Optimal sampling; Semiparametric modelling.

1. Introduction

The emergence of big data area offers statisticians both unprecedented opportunities and challenges. One of the key challenges is that directly applying statistical meth-Corresponding author: zjwang@nankai.edu.cn (Zhaojun Wang) ods to super-large data with conventional computing approaches is prohibitive, which calls for the development of new tools. Recently, statistical analysis and inference in the large-scale dataset have drawn much attention, and some computationally scalable methods have been proposed to reduce the computation and storage effort from various aspects of applications, such as the divide-and-conquer procedures (Battey et al., 2018; Jordan et al., 2019; Zhao et al., 2019), subsampling strategies (Kleiner et al., 2014; Wang et al., 2018) and on-line learning methods (Balakrishnan and Madigan, 2008; Schifano et al., 2016). In most of these works, one usually assumes a parametric model, typically linear or logistic regression models. Therefore, it is necessary to check the misspecification of a given regression model such that the subsequent planning, analysis, and inference can proceed in a creditable way. Lack-of-fit checking for parametric and semiparametric models in a large-scale dataset setting is the focus of this paper.

Suppose Y is the response and $\mathbf{X} = (x_1, \dots, x_p)^{\top} \in \mathbb{R}^p$ is the p-dimensional covariate. We consider a general model in $\overline{\text{Xia}}$ (2009)

$$Y = G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) + \varepsilon, \tag{1.1}$$

where $\mathbf{g} = (g_1, \dots, g_q)^{\mathsf{T}}$ are unknown smooth functions of \mathbf{X} , $G(\cdot)$ is known up to a parameter vector $\boldsymbol{\beta}$, and ε is the random error with $\mathbb{E}(\varepsilon \mid \mathbf{X})=0$. This model includes many parametric and semiparametric models as special cases, such as generalized additive models (Hastie and Tibshirani, 1986), partially linear models (Speckman, 1988), single-index or multi-index models (Hardle et al., 1993; Xia et al., 2002), and varying coefficient models (Hastie and Tibshirani, 1993). Specifically, the generalized additive

models and single-index models admit the forms of $Y = g_1(x_1) + g_2(x_2) + \dots + g_p(x_p) + \varepsilon$ and $Y = g(\mathbf{X}^{\top}\boldsymbol{\beta}) + \varepsilon$, respectively.

The cared model checking problem can be formulated as the following test

$$\mathbb{H}_{0} : \mathbb{E}\left(Y \mid \mathbf{X}\right) = G\left(\mathbf{X}; \boldsymbol{\beta}_{0}, \mathbf{g}_{0}\right), \text{ for some } \boldsymbol{\beta}_{0} \in \Theta, \mathbf{g}_{0} \in \mathcal{G},$$

$$\mathbb{H}_{1} : \mathbb{E}\left(Y \mid \mathbf{X}\right) \neq G\left(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}\right), \text{ for any } \boldsymbol{\beta} \in \Theta, \mathbf{g} \in \mathcal{G},$$

$$(1.2)$$

where $G(\mathbf{X}; \boldsymbol{\beta}_0, \mathbf{g}_0)$ is a pre-specified model with unknown $\boldsymbol{\beta}_0$ and \mathbf{g}_0 , and $\boldsymbol{\Theta}$ and $\boldsymbol{\mathcal{G}}$ are the parameter and function spaces, respectively.

In this paper, we aim to answer the question that "given a limited budget or resources, how can a practitioner optimally use this budget to test (1.2) in a largescale dataset analysis". There are usually two types of limited budgets. On one hand, computing time is a typical budget, which is only able to process a small portion of data. As model checking is very likely one of the most preliminary steps in the data analysis, practitioners would be reluctant to use much computational effort. Several lack-of-fit tests for small and moderate sample sizes have been proposed. Among them, the nonparametric smoothing-based tests and their variants, such as Hardle and Mammen (1993), Zheng (1996) and Fan and Huang (2005), are very popular due to their efficiency and flexibility; see González-Manteiga and Crujeiras (2013) and Guo and Zhu (2017) for some comprehensive reviews. However, the computational complexity and large memory of those methods are typically quadratic in the sample size, which may greatly hamper their applicability to large-scale dataset applications. On the other hand, despite the availability of large-scale datasets, in many applications, collecting responses or labels for all data points is impossible due to measurement constraints

or costs (Wang et al., 2017; Ren et al., 2020), especially at the very beginning of the data processing. As a result, these constraints often require us to select a small subset from a large pool of given design points X and use the limited budget to obtain the corresponding responses Y. For example, in the problem of speech recognition, one may easily get plenty of unlabeled audio data but the accurate labeling of speech utterances is extremely time-consuming and requires trained linguists. Annotation at the word level can take ten times longer than the actual audio (Tur et al., 2005).

When proven statistical methods are no longer applicable due to the two types of limited resources, a natural and appealing method to extract useful information from data is the subsampling method (Kleiner et al., 2014; Ma and Sun, 2015). Many existing pieces of research on subsampling take uniform samples from the full data. However, a nonuniform sampling strategy may achieve better performance. For example, in the estimation problem of linear models, Ma et al. (2015) and Ma and Sun (2015) put forward a so-called algorithmic leveraging with a nonuniform sampling probability to draw a more informative subsample dataset. Some other recent developments in this trend include Wang et al. (2019); Yao and Wang (2019); Yu et al. (2020) and Ai et al. (2021). However, the challenges associated with designing an efficient testing procedure for model checking are not yet well addressed.

In this paper, we propose a new design-adaptive testing procedure for the problem (1.2) when a computation or measurement budget is imposed. The main idea is to select the most informative sample points from the full data, and then construct a computationally tractable test statistic based on the observations of those selected points. We derive an optimal sampling strategy, the Structure-Adaptive-Sampling (SAS), with which the proposed test possesses the asymptotically best power. An initial step is needed to get raw estimations of the quantities involved in the optimal design criterion, and the estimated designs with plug-in estimators are shown to perform equally to the theoretical oracle one from asymptotic viewpoints. The SAS procedure addresses one key question in a general semiparametric framework: how to use the limited resources to implement efficient lack-of-fit tests. Our simulation results clearly demonstrate the superiority of the proposed procedure over existing methods.

The remainder of our paper is structured as follows. In Section 2, we present the construction of the optimal sampling designs along with a detailed discussion on asymptotic justifications. Some practical guidelines are given in Section 3. Numerical studies and a real-data example are conducted in Section 4. Section 5 concludes the paper, and theoretical proofs are delineated in the Appendix.

2. Methodology

Assume that there are total N available data points or observable subjects $\mathcal{X} = \{\mathbf{X}_i^a\}_{j=1}^N \in \mathbb{R}^p$. Given a measurement constraint, only n samples $\mathcal{S} = \{\mathbf{X}_i, Y_i\}_{i=1}^n$ can be obtained, or similarly the computational budget only allows us to deal with one dataset with size n, where Y_i is the response and $n \ll N$. For the dataset \mathcal{S} , it is assumed that we independently sample \mathbf{X}_i from \mathcal{X} with replacement and then observe its corresponding response Y_i . We start with the test construction on \mathcal{S} given the full data.

2.1 Test construction

Denote that the residual $\epsilon = Y - G(\mathbf{X}; \boldsymbol{\beta}_0, \mathbf{g}_0)$, and then test problem (1.2) amounts to assess whether $\mathbb{E}(\epsilon \mid \mathbf{X}) = \mathbb{E}\{Y - G(\mathbf{X}; \boldsymbol{\beta}_0, \mathbf{g}_0) \mid \mathbf{X}\}$ is equal to zero or not. A standard way is to construct a test statistic based on estimated $\mathbb{E}(\epsilon \mid \mathbf{X})$ with fitted residuals under the null hypothesis. See, for example, Hardle and Mammen (1993); Stute et al. (1998); Dette (1999); Fan et al. (2001). However, due to the difficulty in nonparametric estimation of the function $\mathbb{E}(\epsilon \mid \mathbf{X})$ or $\mathbb{E}(Y \mid \mathbf{X})$ when p > 2, the efficiency of those methods drops rapidly as the dimension p of the covariates increases.

To this end, we consider a structured alternative model as $\mathbb{E}\left(\epsilon \mid \mathbf{X}\right) = M\left(\boldsymbol{\theta}^{\top}\mathbf{X}\right)$, where $\boldsymbol{\theta} \in \mathbb{R}^{p}$ is one projection direction with $\|\boldsymbol{\theta}\|_{2} = 1$ and $M(\cdot)$ is an unknown smooth function. Thus, the alternative $\mathbb{E}\left(\epsilon \mid \mathbf{X}\right) \neq 0$ is equivalent to $\mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right) \neq 0$. It is also clear that $\mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right) = 0$ due to $\mathbb{E}\left(\epsilon \mid \mathbf{X}\right) = 0$ under the null hypothesis. Then, test problem (1.2) can be formulated to the following one

$$\mathbb{H}_0 : \mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top} \mathbf{X}\right) = 0 \text{ versus } \mathbb{H}_1 : \mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top} \mathbf{X}\right) = M(\boldsymbol{\theta}^{\top} \mathbf{X}) \neq 0,$$
 (2.1)

where $M(\boldsymbol{\theta}^{\top}\mathbf{X}) \neq 0$ for some $\boldsymbol{\theta}^{\top}\mathbf{X} \in \Omega \subset \mathbb{R}$. Intuitively, $\mathbb{E}\left\{\epsilon\mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right)\right\}$ is a good choice to measure the derivation between \mathbb{H}_0 and \mathbb{H}_1 . Under \mathbb{H}_0 , $\mathbb{E}\left\{\epsilon\mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right)\right\} = 0$, while under \mathbb{H}_1 , $\mathbb{E}\left\{\epsilon\mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right)\right\} = \mathbb{E}\left\{\mathbb{E}^2\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right)\right\} > 0$. Then, our test can be built by $\mathbb{E}\left\{\epsilon\mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right)\right\}$, which is a popular quantity in the context of model specification test (Zheng, 1996; Guerre and Lavergne, 2005).

Given the projection direction $\boldsymbol{\theta}$, a nonparametric estimation of $\mathbb{E}\left\{\epsilon\mathbb{E}\left(\epsilon\mid\boldsymbol{\theta}^{\top}\mathbf{X}\right)\right\}$

based on $S = \{\mathbf{X}_i, Y_i\}_{i=1}^n$ is

$$V_{f} = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i}^{n} \frac{\widehat{\epsilon}_{i} \widehat{\epsilon}_{j} K_{h} (\omega_{i} - \omega_{j})}{\sqrt{f(\omega_{i})} \sqrt{f(\omega_{j})}},$$

where $\omega_i = \boldsymbol{\theta}^{\top} \mathbf{X}_i$ is followed by a density $f(\cdot)$, $\widehat{\epsilon}_i = Y_i - G(\mathbf{X}_i; \widehat{\boldsymbol{\beta}}, \widehat{\mathbf{g}})$ are the fitted residuals, $\widehat{\boldsymbol{\beta}}$ and $\widehat{\mathbf{g}}$ are respectively estimates of $\boldsymbol{\beta}$ and \mathbf{g} , and $K_h(\cdot) = K(\cdot/h)/h$ denotes a one-dimensional kernel function with a bandwidth h. For notation simplicity, we only emphasize the dependence of test statistic V_f on density $f(\cdot)$. Under some mild conditions, it can be shown that V_f is asymptotically normal under \mathbb{H}_0 (Zheng, 1996; Guerre and Lavergne, 2005), that is

$$T_f := nh^{1/2}V_f/\sigma_V \xrightarrow{\mathcal{L}} \mathcal{N}(0,1), \tag{2.2}$$

where $\sigma_V^2 = 2\sigma^4 |\Omega| \int K^2(u) du$ is the asymptotic variance under null, and $|\Omega|$ is the cardinality of Ω and $\sigma^2 = \mathbb{E}(\epsilon^2 \mid \boldsymbol{\theta}^\top \mathbf{X})$. The arrow $\stackrel{\mathcal{L}}{\longrightarrow}$ should be understood as convergence in distribution. A large value of T_f would lead to rejection of the null.

The use of $\boldsymbol{\theta}$ plays an important role of dimension reduction in the test statistic V_f . In practice, $\boldsymbol{\theta}$ can be either user-specified or estimated via some dimension reduction techniques. Detailed discussions on the determination of $\boldsymbol{\theta}$ can be found in Section 3. Similar projection-based tests have been proposed in the literature. For example, Fan and Huang (2001) reduced the dimension based on \mathbf{X} alone; Xia (2009) proposed to project the fitted residuals along a direction that adapts to the systematic departure of the residuals from the hypothetical pattern with cross-validation in the single-index model. A more related work is Guo et al. (2016) in which they checked the single-index models based on a joint estimation of the dimension reduction matrix. Other works include Zhu et al. (2017) and Tan et al. (2018).

2.2 Optimal sampling strategy

To implement the proposed method, it is important to specify the sampling density $f(\omega)$. The conventional choice is to let $f(\omega)$ as uniform distribution (Stute and Zhu, 2002; Ferragut and Laska, 2012; Guo et al., 2016), which corresponds to the simple uniform sampling. However, the uniform sampling may not be *optimal* so that the informative subsamples are not selected. Under certain local alternatives, the asymptotic distribution of T_f is also normal but with a positive mean, which depends on $f(\omega)$. It implies that if an appropriate sampling density $f(\omega)$ can be chosen, the power of this test will be maximized.

Next, we will provide a result that sheds lights on how to determine the optimal sampling density. Firstly, we need the following to facilitate the derivation. Let $(\boldsymbol{\beta}^*, \mathbf{g}^*) = \arg\min_{(\boldsymbol{\beta}, \mathbf{g}) \in \Theta \otimes \mathcal{G}} \mathbb{E} \{Y - G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g})\}^2$. Under \mathbb{H}_0 , it is clear that $(\boldsymbol{\beta}^*, \mathbf{g}^*) = (\boldsymbol{\beta}_0, \mathbf{g}_0)$.

Assumption 1. (Moments condition) For $\kappa = 4 + \gamma$ with small enough $\gamma > 0$, $\mathbb{E}(\varepsilon^{\kappa} \mid \mathbf{X}) \leq C_1 < \infty$, where $C_1 > 0$ is a fixed constant.

Assumption 2. (Density function) The density function of ω , $f(\omega)$, is continuous on the compact support $\omega \in \Omega$, satisfying $0 < \inf_{\omega \in \Omega} f(\omega) \le \sup_{\omega \in \Omega} f(\omega) < \infty$.

Assumption 3. (Kernel function) $K(\cdot)$ is a continuous, nonnegative, bounded, and symmetric kernel function with a bounded first order derivative.

Assumption 4. (Model) The semiparametric model can be approximated by first order

Taylor expansion

$$G(\mathbf{X}; \widehat{\boldsymbol{\beta}}, \widehat{\mathbf{g}}) = G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) + \nabla G_{\boldsymbol{\beta}}^{\top}(\mathbf{X}; \widetilde{\boldsymbol{\beta}}, \widetilde{\mathbf{g}}) (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + \nabla G_{\mathbf{g}}^{\top}(\mathbf{X}; \widetilde{\boldsymbol{\beta}}, \widetilde{\mathbf{g}}) (\widehat{\mathbf{g}} - \mathbf{g}),$$

where $\nabla G_{\mathbf{g}}(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) = \left\{ \nabla G_{g_1}(\mathbf{X}; \boldsymbol{\beta}, \mathbf{z}), \dots, \nabla G_{g_q}(\mathbf{X}; \boldsymbol{\beta}, \mathbf{z}) \right\}^{\top}$ with $\nabla G_{g_k}(\mathbf{X}; \boldsymbol{\beta}, \mathbf{z}) = \partial G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{z})/\partial z_k$. Here $\widetilde{g}_k(\nu)$ lies between $g_k(\nu)$ and $\widehat{g}_k(\nu)$, $k = 1, \dots, q$, and $\widetilde{\boldsymbol{\beta}}$ lies between $\boldsymbol{\beta}$ and $\widehat{\boldsymbol{\beta}}$. And $\nabla G_{\boldsymbol{\beta}}(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) = \partial G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g})/\partial \boldsymbol{\beta}$. Furthermore, $\nabla G_{\boldsymbol{\beta}}$ and $\nabla G_{\mathbf{g}}$ are Lipschitz continuous with $0 < \max_{\boldsymbol{\beta} \in \Theta, \mathbf{g} \in \mathcal{G}} \left(\mathbb{E} \left\{ \nabla G_{\boldsymbol{\beta}}(\mathbf{X}_i; \boldsymbol{\beta}, g_k) \right\}^2, \mathbb{E} \left\{ \nabla G_{\mathbf{g}}(\mathbf{X}_i; \boldsymbol{\beta}, g_k) \right\}^2 \right) \leq C_2 < \infty, i = 1, \dots, n, k = 1, \dots, q \text{ with some positive constant } C_2.$

Assumption 5. (Asymptotic representation) Suppose $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\| = O_p(n^{-1/2})$. Assume that all the link functions g_1, \ldots, g_q own a common compact support Γ , and their estimators admit the following asymptotic expansions

$$\sup_{\nu \in \Gamma} \left| \widehat{g}_k(\nu) - g_k^*(\nu) - R_k(\nu)b^2 - n^{-1}H_k(\nu) \sum_{i=1}^n \phi_k(\mathbf{X}_i) Q_b(u_{ki} - \nu) \varepsilon_i \right| = O_p(n^{-1/2}),$$

for k = 1, ..., q, where u_{ki} is a measurable function of \mathbf{X}_i , $R(\cdot)$, $H(\cdot)$, and $\phi(\cdot)$ are bounded continuous functions. For some positive integer $r \geq 2$, the rth derivative of $g_k(\cdot)$ is bounded. $Q_b(\cdot)$ is a bounded, symmetric and r order continuously differentiable kernel function with smoothing parameter b, satisfying $\int Q_b(u) du = 1$, $\int u^i Q_b(u) du = 0$, and $\int u^r Q_b(u) du \neq 0$ for $0 \leq i < r$, $b \to 0$ and $nb \to \infty$ as $n \to \infty$.

Assumption 6. (Bandwidth) The testing bandwidth h satisfies $h \to 0$, $nh \to \infty$, and $nh^{1/2}b^{2r} \to 0$ as $n \to \infty$ with positive integer $r \ge 2$.

Remark 1. Assumptions $\boxed{1}$ — $\boxed{3}$ are quite standard in kernel-based methods, though some of them may not be the weakest possible. For instance, we only require $f(\cdot)$

to be Lipschitz continuous and bounded away from zero if some other conditions are imposed. Assumption 4 is a regularity condition on the semiparametric model. This is reasonable because we cannot expect that our procedure would work well if the $G(\cdot)$ is not in a regular form. The formulation is quite mild and can be satisfied by all the commonly used semiparametric models mentioned before. Assumption 5 sets theoretical requirements for the estimates of the model, and is fulfilled by most semiparametric estimation methods and models, including the partially linear model (Speckman, 1988), the additive model (Horowitz and Mammen, 2004), the varying coefficient model (Fan and Zhang, 1999) and the single-index model (Ichimura, 1993) under the condition that the nonparametric part \mathbf{g} is twice continuously differentiable function on Ω ; see also $\overline{\text{Xia}}$ (2009) for similar discussions. Assumption 6 gives the bandwidth requirements for implementing the test V_f with asymptotic normal calibration. The optimal rate of nonparametric estimation $n^{-1/5}$ appears to be not allowed if we set b=h for simplicity and consider the case r=2. This is very common in the problem of nonparametric specification tests; certain under-smoothing is necessary (see Fan et al., 2001).

Consider the local alternative model as $Y = G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) + \delta_n l(\omega) + \varepsilon$ and the corresponding local alternatives become

$$\mathbb{H}_{1n}: M(\omega) = \delta_n l(\omega) \quad \text{for} \quad \omega \in \Omega,$$
 (2.3)

where $M(\omega) = \mathbb{E}\left(\epsilon \mid \boldsymbol{\theta}^{\top}\mathbf{X} = \omega\right)$ and $\delta_n \to 0$ as $n \to \infty$, $l(\omega)$ is continuously differentiable on Ω satisfying $\mathbb{E}\left\{l^2(\omega)\right\} < \infty$ and bounded away from zero almost everywhere. Under (2.3), $\delta_n l(\omega)$ characterizes the model difference from the null to alternative, which goes to zero as $n \to \infty$. Then, we have the following result.

Theorem 1. Suppose Assumptions 1—6 hold. Under the local alternatives (2.3) with $\delta_n = (nh^{1/2})^{-1/2}$, the test based on T_f reaches its asymptotic best power if the sampling density $f(\omega) \propto M^2(\omega)$.

Under the local alternative [2.3], $\mathbb{E}(T_f) = \mathbb{E}_f \{l^2(\omega)\}/\sigma_V$ and the asymptotic variance σ_V^2 is irrelevant with $f(\omega)$ and $l(\omega)$, so an explicit power expression, depending on $\mathbb{E}_f \{l^2(\omega)\}$, can be obtained. Theorem [1] enlightens us to construct a locally most powerful test by choosing the sampling distribution $f(\omega)$ as a linear function of $M^2(\omega)$. A similar optimization technique, Cauchy's inequality, is adopted in Wang et al. (2018) to derive the nonuniform sampling strategy under some optimality criterion, but it is an estimation problem and its focus is to estimate the sampling probabilities when the whole data is available which is quite different with our testing procedure. In our work, we only need to sample the data point \mathbf{X} with sampling density $f(\cdot)$ and then observe its corresponding response Y. Since the sampling density is related to the underlying model structure, we term this strategy as the *Structure-Adaptive-Sampling* (SAS) procedure.

2.3 Estimation of optimal density $f(\cdot)$

The optimal sampling density depending on $M(\omega)$ contains the unknown function $l(\omega)$ in Theorem \mathbb{I} , so an exact $f(\omega)$ is inapplicable directly for practical implementation. Hence, it is necessary to obtain raw but informative estimates by a pilot study. Then an approximately optimal sampling plan can be achieved. Assume we have the dataset $S_0 = \{\mathbf{X}_{0i}, Y_{0i}\}_{i=1}^{n_0}$ for pilot study, where $n_0 \leq n$ and $\{\mathbf{X}_{0i}\}_{i=1}^{n_0}$ are uniformly

sampled from \mathcal{X} . Given $\boldsymbol{\theta}$, let $\{\widehat{\epsilon}_{0i}\}_{i=1}^{n_0}$ be the fitted residuals based on \mathcal{S}_0 , then a consistent estimator of $M(\omega)$ is

$$\widehat{M}(\omega) = \frac{\sum_{i=1}^{n_0} K_{h_f} (\omega_{0i} - \omega) \widehat{\epsilon}_{0i}}{\sum_{i=1}^{n_0} K_{h_f} (\omega_{0i} - \omega)},$$
(2.4)

where $\omega_{0i} = \boldsymbol{\theta}^{\top} \mathbf{X}_{0i}$ and h_f is a pre-specified bandwidth which satisfies $h_f \to 0$ and $n_0 h_f / \log n_0 \to \infty$ as $n_0 \to \infty$. By Theorem 1, the optimal distribution $f(\omega)$ can be estimated by

$$\widehat{f}(\omega) = \frac{\max\{\xi_{n_0}, \widehat{M}^2(\omega)\}}{\int \max\{\xi_{n_0}, \widehat{M}^2(\omega)\} d\omega},$$
(2.5)

where ξ_{n_0} is fixed as $O\{(\log n_0)^c/(n_0h_f)\}$ for c>1. The use of ξ_{n_0} in (2.5) is to ensure that the estimated density is bounded away from zero, especially under \mathbb{H}_0 where $M(\omega) = 0$ for any $\omega \in \Omega$. By our theoretical results given in the Appendix, $\sup_{\omega} |M(\omega)| = o_p(\xi_{n_0}^{1/2})$ under \mathbb{H}_0 , and further $\widehat{f}(\omega)$ is degenerated to a uniform distribution with probability tending to one, i.e., $\widehat{f}(\omega) = 1/|\Omega|$.

The next result shows that the effect of replacing $M(\omega)$ by appropriate estimators can be asymptotically negligible and the efficiency of the locally most powerful test can still be achieved under some mild conditions.

Theorem 2. Assume Assumptions 1—6 and $(n_0/n)(h_f^2/h)^{1/2}/(\log n_0)^c \to \infty$ all hold. Under the local alternatives \mathbb{H}_{1n} (2.3), the power function of our SAS test with $\widehat{f}(\omega)$ can be approximated by

$$\Pi_f = \Phi\left(\left\{\int l^4(\omega)d\omega/\int l^2(\omega)d\omega\right\}/\sigma_V - z_\alpha\right),$$

where $\Phi(\cdot)$ denotes the cumulative distribution function of standard normal, and z_{α} is the corresponding upper- α quantile.

Note that $\mathbb{E}_f\{l^2(\omega)\}$ is the alternative mean of test statistics V_f and satisfies $\mathbb{E}_f\{l^2(\omega)\} \leq \int l^4(\omega) d\omega / \int l^2(\omega) d\omega$. Thus, using the estimated density $\widehat{f}(\omega)$ in (2.5) yields a more powerful test.

Furthermore, a more compelling result is that the optimal sampling strategy could be more prominent when the signal function $M(\omega)$ exhibits a sparse pattern. Consider the following sequence of "singular" local alternatives

$$\mathbb{H}'_{1n}: M(\omega) = \delta'_n l(\omega) \text{ for } \omega \in \Omega_n,$$
 (2.6)

where $\Omega_n \subset \Omega$ satisfying $|\Omega_n| \approx a_n$ with $a_n \to 0$ being a deterministic sequence and $l(\omega)$ is bounded away from zero on Ω_n almost everywhere. The main feature of these "singular" local alternatives is that they have narrow spikes and change rapidly as the sample size n increases. In other words, the alternatives in (2.6) can be regarded as sparse/high-frequency alternatives, while the alternatives in (2.3) as dense/low-frequency alternatives.

Corollary 1. Consider the "singular" local alternative [2.6] with $\delta'_n = (nh^{1/2}a_n^{1/2})^{-1/2}$, $a_n/h \to \infty$ and $\inf_{\omega \in \Omega_n} l(\omega) \ge l_{\min} > 0$. Suppose the conditions in Corollary [2] hold. If $(n_0/n)(h_f^2/h)^{1/2}/\left\{a_n^{1/2}(\log n_0)^c\right\} \to \infty$, the asymptotic power of our sampling test with $\widehat{f}(\omega)$ is not small than $\Phi(\mu(f,\Omega_n)/\sigma_V(f,\Omega_n)-z_\alpha)$, where $\mu(f,\Omega_n)=|\Omega_n|^{-1/2}l_{\min}^2$ and $\sigma_V^2(f,\Omega_n)=2\sigma^4|\Omega_n|\int K^2(u)du$.

The condition $a_n/h \to \infty$ in Corollary 1 ensures that our test could work well if the sparse signal size of Ω_n goes to zero slower than h as n increases. The advantage of the proposed test under (2.6) is that the conditions imposed on estimating bandwidth h_f are more relaxed in Corollary 1. That is if $a_n = h^{1/2}$, the optimal rate of nonparametric

bandwidth $h_f = O(n_0^{-1/5})$ can be allowed as long as the testing bandwidth h satisfies $nh^{3/4}/n_0^{4/5} \to 0$.

2.4 The SAS-bsed testing procedure

Hence, our procedure for model checking is summarized as follows.

Algorithm 1 Model checking via structure-adaptive-sampling

Step 1 (Initialization) Specify $K(\cdot)$, n_0 , n, h_f , h and θ ;

Step 2 (Pilot study) Estimate model (1.1) based on $S_0 = \{\mathbf{X}_{0i}, Y_{0i}\}_{i=1}^{n_0}$ and compute the fitted residuals $\{\hat{\epsilon}_{0i}\}_{i=1}^{n_0}$, where \mathbf{X}_{0i} 's are uniformly sampled from \mathcal{X} ;

Step 3 (Sampling) Obtain $\widehat{f}(\omega)$ by (2.5) with $\widehat{M}(\omega)$ in (2.4) based on \mathcal{S}_0 and $\{\widehat{\epsilon}_{0i}\}_{i=1}^{n_0}$; sample n data points \mathbf{X}_i 's with $\widehat{f}(\omega)$ from \mathcal{X} and get the corresponding Y_i 's as $\mathcal{S} = \{\mathbf{X}_i, Y_i\}_{i=1}^n$;

Step 4 (Test) Compute $\{\widehat{\epsilon_i}\}_{i=1}^n$ based on \mathcal{S} , and then build $T_{\widehat{f}}$ by (2.2); further, reject \mathbb{H}_0 if $T_{\widehat{f}} > z_{\alpha}$.

To generate n samples at Step 3, we could estimate the sampling probabilities by (2.5) for N possible data points in \mathcal{X} and then draw n sampling points $\{\mathbf{X}_i\}_{i=1}^n$ by a multinomial distribution. This sampling procedure is fast to implement with the computational complexity $O(Nn_0h_f)$. Note that the computational complexity for the estimation of $G(\cdot)$ is generally linear in the sample size, thus the computing time of the pilot study is $O(n_0)$. In this way, the computation of our SAS testing procedure is $O(n_0 + Nn_0h_f + n^2h)$, where $O(n^2h)$ is the complexity for computing $T_{\widehat{f}}$ in Step 4.

We use a numerical example to demonstrate the performance of our SAS test. We

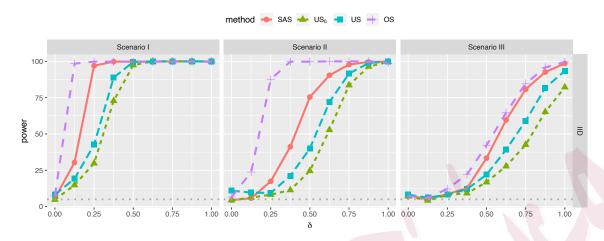


Figure 1: Empirical sizes and powers (%) for the tests with different sampling methods under Scenarios I–III when errors follow $\mathcal{N}(0,1)$ independently and the design points \mathbf{X} are from $\mathcal{N}(0,\mathbf{I}_p)$. The gray dotted line is the significance level $\alpha=0.05$.

sample n=1000 data points and $n_0=300$ pilot samples from \mathcal{X} with size $N=10^5$ and consider three alternative models. See Scenarios I–III in Section 4 for details. Figure 1 compares the power curves of our proposed procedure, SAS, with three other sampling-based tests. They differ only in their sampling strategies. To be specific, OS uses the so-called "oracle" density $M^2(\omega)/\int M^2(\omega)d\omega$ as if $M(\omega)$ is known, and the US₀ and US use uniform distribution to sample n and $n+n_0$ for building T_f respectively. The pilot study (if needed) for all methods is based on n_0 observations. The improvement of our adaptive-sampling-based test over US₀ and US test is clear. The OS approach has superior performance as expected, but the difference between SAS and OS becomes smaller than US₀ and the US as signal δ increases.

In the foregoing discussion, we assume the alternative model as $\mathbb{E}\left(\epsilon \mid \mathbf{X}\right) = M(\boldsymbol{\theta}^{\top}\mathbf{X}) \neq 0$ for any vector $\boldsymbol{\theta}$. In fact, we cannot expect this model holds exactly, and a more re-

alistic assumption is that $\mathbb{E}(\epsilon \mid \mathbf{X}) = M_{\mathbf{B}}(\mathbf{B}^{\top}\mathbf{X})$, where \mathbf{B} is a $p \times d$ matrix with unknown $d \geq 1$. In such a situation, we may work with a misspecified model. However, as long as the alternative satisfies $\mathbb{E}\{M_{\mathbf{B}}(\mathbf{B}^{\top}\mathbf{X}) \mid \boldsymbol{\theta}^{\top}\mathbf{X}\} \neq 0$, the proposed test would still work well and its power function can be verified by $\Pi_{M_{\mathbf{B}}} = \Phi(\mu(f, \mathbf{B}, \boldsymbol{\theta}) / \sigma_V(\mathbf{B}, \boldsymbol{\theta}) - z_{\alpha})$, where $\sigma_V^2(\mathbf{B}, \boldsymbol{\theta}) = 2\sigma^4 |\Psi| \int K_d^2(s) ds$ with $\mathbf{B}^{\top}\mathbf{X} \in \Psi$ and a d-dimensional kernel function $K_d(\cdot)$, and $\mu(f, \mathbf{B}, \boldsymbol{\theta}) = \mathbb{E}_f\left[\mathbb{E}\{M_{\mathbf{B}}(\mathbf{B}^{\top}\mathbf{X})\} \mid \boldsymbol{\theta}^{\top}\mathbf{X}\right]^2$. Note that if $\mathbb{E}\{\mathbb{E}(\epsilon \mid \mathbf{X}) \mid \boldsymbol{\theta}^{\top}\mathbf{X}\}$ is a function of $\boldsymbol{\theta}^{\top}\mathbf{X}$, say d = 1, then $\Pi_{M_{\mathbf{B}}}$ reduces to the one given in Theorem \mathbb{Z} .

3. Practical guidelines

In this section, several practical issues on implementing the SAS procedure are discussed, including the choices of projection direction and bandwidths as well as the determination of n_0 .

3.1 Determination of θ

One key in the implementation of the SAS procedure is the selection of the projection direction $\boldsymbol{\theta}$. Actually, one can assign a fixed $\boldsymbol{\theta}$ based on practical requirements or estimate it via some dimension reduction techniques. For example, we want to check whether there might be a nonlinear relationship between the response Y and covariate x_1 when the null hypothesis is a linear function. In this case, we can directly set $\boldsymbol{\theta} = (1, 0, \dots, 0)^{\top}$.

In general, we can obtain a reliable estimation of θ in the pilot study with some techniques on sufficient dimension reduction (SDR). To identify the dimension reduc-

tion subspace, the literature contains many proposals, such as classical sliced inverse regression (SIR, Li, 1991), sliced average variance estimation (SAVE, Cook and Weisberg, 1991), directional regression (DR, Li and Wang, 2007) and likelihood acquired directions (Cook and Forzani, 2009). Specifically, Xia et al. (2002) proposed the minimum average variance estimate (MAVE) to estimate the reduction space with fewer regularity conditions on the covariates X.

Next, we demonstrate the asymptotic properties of the test statistics $T_{\widehat{f}}$ when $\boldsymbol{\theta}$ is estimated in the pilot study (Step 2 in Algorithm $\boxed{1}$).

Assumption 7. (Projection direction) The estimated projection direction with a sample of size n_0 satisfies $\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta} = O_p(n_0^{-1/2})$ as $n_0 \to \infty$, where

$$\boldsymbol{\theta} = \operatorname*{arg\,min}_{\boldsymbol{\alpha}: \|\boldsymbol{\alpha}\|_2 = 1} \mathbb{E} \left\{ \epsilon - \mathbb{E}(\epsilon \mid \boldsymbol{\alpha}^\top \mathbf{X} = \omega) \right\}^2.$$

Assumption 7 is very mild and typically holds for most of the SDR methods.

Theorem 3. Suppose Assumptions $\boxed{1-\widehat{\gamma}}$ and $n_0h^{3/2} \to \infty$ hold. The variance σ_V^2 can be consistently estimated by $\widehat{\sigma}_V^2 = \left\{n(n-1)\right\}^{-1} 2h \sum_{i=1}^n \sum_{j\neq i}^n \widehat{\epsilon}_i^2 \widehat{\epsilon}_j^2 K_h^2 \left(\widehat{\omega}_i - \widehat{\omega}_j\right) \left\{\widehat{f}(\widehat{\omega}_i)\widehat{f}(\widehat{\omega}_j)\right\}^{-1}$. We have

- (i) Under the null hypothesis \mathbb{H}_0 (2.1), $T_{\widehat{f}} \xrightarrow{\mathcal{L}} \mathcal{N}(0,1)$;
- (ii) Under the local alternative \mathbb{H}_{1n} (2.3) with $(n_0/n)(h_f^2/h)^{1/2}/(\log n_0)^c \to \infty$ and $\delta_n = (nh^{1/2})^{-1/2}, \ T_{\widehat{f}} \{\int l^4(\omega) d\omega / \int l^2(\omega) d\omega \}/\widehat{\sigma}_V \xrightarrow{\mathcal{L}} \mathcal{N}(0,1).$

Theorem 3 can be viewed as a counterpart of the asymptotic normality result given in Zheng (1996) and Guerre and Lavergne (2005) as shown in (2.2), by generalizing a

parametric null specification model to a much more generic semiparametric one (1.1). Due to the involvement of both the estimations of θ and $G(\cdot)$, the technical details of our theory are not straightforward and cannot be obtained from those existing works.

Theorem 3 implies that our SAS procedure could achieve the best power with $\hat{\boldsymbol{\theta}}$ from asymptotic viewpoints and the results in Theorem 2 still can achieve. It is further verified by the numerical comparisons over a wide range of values of $\boldsymbol{\theta}$ in the Supplementary Material Table 51. In this paper, we use the MAVE method (Xia et al., 2002) to estimate $\boldsymbol{\theta}$ in the pilot study due to its easy implementation with the R package MAVE.

3.2 Bandwidth selection

Like many other smoothing-based tests, the performance of the SAS test possibly depends on the bandwidth h in the test statistic $T_{\widehat{f}}$ and the h_f in density estimator (2.4). By Corollary 1, the optimal h_f for estimator $\widehat{M}(\omega)$ can achieve the order $O(n_0^{-1/5})$. Thus, we take the empirical bandwidth formula $h_f = 0.5 \text{sd}(\omega_0) n_0^{-1/5}$ as a rule of thumb, where $\text{sd}(\omega_0)$ denotes the sample standard deviation of $\omega_0 = \{\boldsymbol{\theta}^{\top} \mathbf{X}_{0i}\}_{i=1}^{n_0}$.

Different from the estimating bandwidth h_f , it is widely known that the selection of bandwidth h for optimal testing power is an open problem (Hart, 1997; Stute et al., 2005). Asymptotically, a range of bandwidths that satisfy Assumption 6 will retain the consistency of the test, while a larger bandwidth generally results in a better power in the view of Corollary 2. However, in practice, the condition $h \to 0$ will restrict h not too large, and the condition $nh^{3/4}/n_0^{4/5} \to 0$ will ensure that h cannot be too small. Based on our numerical results, the observed significance changes only mildly over a

wide range of value of h, and we recommend $h = \{h_f(n_0/n)^{1/5}\}^{2+\eta}$ with some $\eta > 0$ so that the condition $(n_0/n)(h_f^2/h)^{1/2}/\{a_n^{1/2}(\log n_0)^c\} \to \infty$ in Corollary 1 is roughly valid. This choice works well for a wide range of models and pilot sample sizes as shown in Section 4.

3.3 Choices of sample sizes

In practice, we could uniformly sample n_0 data points as S_0 for the pilot study. Note that there is a trade-off in the selection of n_0 between estimation efficiency and computational complexity. Intuitively, a larger n_0 would attain more accurate estimations of the density function $\hat{f}(\cdot)$ and projection direction θ (if needed), but involves more computational burden. To satisfy our theoretical requirements, we could roughly consider $n_0 = \lfloor n^{3/5} (\log n)^{c_0} \rfloor$ with some $c_0 > 0$. Our simulations appear that reliable estimations could be obtained and the performance of SAS is not affected too much as long as $n_0 \geq 200$ in the pilot study, which seems to be acceptable for a large-scale dataset.

It is also worth noting that better performance would be expected when n is larger. How large sample size n is allowed depends on practitioners' resource constraints, such as computational power, measurement costs, and processing time.

4. Numerical studies

In this section, we examine the performance of our proposed SAS procedure via some Monte Carlo simulation studies and a real-data example.

4.1 Simulation studies

The data are generated by the model $Y = G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) + \varepsilon$, where ε is distributed from $\mathcal{N}(0,1)$ and the covariates \mathbf{X} are independently from the whole space \mathcal{X} . Consider the whole data points \mathcal{X} with $N = |\mathcal{X}| = 10^5$ are i.i.d from $\mathcal{N}(0, \mathbf{\Sigma})$. Two classes of $\mathbf{\Sigma}$ are explored: one is the identity matrix \mathbf{I}_p , and the other has the components $(\mathbf{\Sigma})_{ij} = 0.5^{|i-j|}$ for $i, j = 1, \ldots, p$. We denote these two as "IID" and "COR" cases respectively. The sample size n is fixed as 1000. Three scenarios for $G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g})$ are considered:

- Scenario I (Linear Model): $G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) = \boldsymbol{\beta}^{\top} \mathbf{X} + \delta \cdot 0.4 |\boldsymbol{\theta}^{\top} \mathbf{X}|^3$ under p = 4, $\boldsymbol{\beta} = (1, 1, -1, -1)^{\top}/2$ and $\boldsymbol{\theta} = (1, -1, 0, 0)^{\top}/\sqrt{2}$;
- Scenario II (Single-index Model): $G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) = 2 \exp(-\boldsymbol{\beta}^{\top} \mathbf{X}) + \delta \cdot 0.4 (\boldsymbol{\theta}^{\top} \mathbf{X})^2$ under p = 6, $\boldsymbol{\beta} = (1, -1, 0, 0, 0, 0)^{\top} / \sqrt{2}$ and $\boldsymbol{\theta} = (0, 0, 1, -1, 0, 0)^{\top} / \sqrt{2}$;
- Scenario III (Multi-index Model): $G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) = (\boldsymbol{\beta}_1^{\top} \mathbf{X})^2 + (\boldsymbol{\beta}_2^{\top} \mathbf{X})^2 + \delta \cdot 0.8 \exp\left(-0.4\boldsymbol{\theta}^{\top} \mathbf{X}\right)$ under p = 4, $\boldsymbol{\beta}_1 = (1, 0, 0, 0)^{\top}$, $\boldsymbol{\beta}_2 = (0, 1, 0, 0)^{\top}$ and $\boldsymbol{\theta} = (0, 0, 1, 0)^{\top}$.

where $\delta \geq 0$, and $\delta = 0$ corresponds the null hypothesis in (1.2). All the simulation results are based on 500 replications.

We fix the target significant level α as 0.05 and evaluate the performance of SAS procedure via comparing empirical sizes under the null and powers under the alternatives. We discuss the choices of kernel functions for the nonparametric test statistics in the Supplementary Material Table S2 and find that the effect of kernel functions can be ignorable. The Epanechnikov kernel function is applied here. As discussed in

Section 3 we consider the empirical bandwidth formula $h_f = 0.5 \text{sd}(\omega_0) n_0^{-1/5}$ in the density estimator (2.5) and $h = \{h_f(n_0/n)^{1/5}\}^{2+\eta}$ in the test statistics $T_{\widehat{f}}$ for some $\eta > 0$, where $\text{sd}(\omega_0)$ is the sample standard deviation of $\omega_0 = \{\boldsymbol{\theta}^{\top} \mathbf{X}_{0i}\}_{i=1}^{n_0}$. Table 1 reports the empirical sizes and powers of SAS procedure with different bandwidths when Σ is from the "IID" case. We observe that three different values of $\eta \in (0, 0.2]$ present similar results: the empirical sizes and powers are not significantly different across all the settings. Meanwhile, our SAS procedure is not affected too much under the null and presents reliable power under the alternatives when the sample size n_0 in the pilot study is larger than 200. Hence, $\eta = 0.1$ and $n_0 = 300$ are used in the rest of simulations.

We next compare our SAS procedure with several benchmarks. Besides the methods with uniform sampling strategy as mentioned before, i.e., US₀ and US, we also consider two other existing methods. The first one is from Guo et al. (2016), in which they proposed a dimension reduction model-adaptive local smoothing test for parametric single-index models. We term it as Guo for simplicity. The second one is a global testing approach from Stute and Zhu (2002), named as SZ here, in which they developed a dimension reduction test and approximated the distribution of the test statistics based on a certain empirical process. To make a fair comparison, we apply Guo and SZ on a subset with $n + n_0$ samples where the covariates \mathbf{X} are uniform sampling from \mathcal{X} and use wild bootstrap of 500 times to mimic their critical values.

Figure 2 displays the comparisons of empirical sizes and powers among SAS, US₀ and US. Our SAS outperforms US₀ and US uniformly across all the settings. This is not

Table 1: Empirical sizes and powers (%) of SAS procedure with different bandwidths under Scenarios I–III when Σ is from the "IID" case.

		Scenario I			Scenario II				Scenario III		
n_0	δ	0.00	0.25	0.50	0.00	0.25	0.50	(0.00	0.25	0.50
200	0.05	5.8	91.0	100.0	4.2	9.2	58.2		5.0	5.4	26.8
	0.1	5.6	90.8	100.0	4.4	9.2	57.8		4.8	5.4	25.4
	0.2	6.6	90.0	99.6	4.6	9.0	57.2		4.8	5.4	23.4
300	0.05	6.6	97.0	100.0	4.2	16.6	75.2		4.4	9.0	32.2
	0.1	5.8	97.0	99.8	4.4	17.4	75.4		5.2	9.0	31.8
	0.2	5.8	95.6	99.8	4.2	16.8	74.2		4.6	9.0	30.2
400	0.05	7.0	98.0	100.0	4.6	19.8	85.8		4.6	8.8	38.6
	0.1	7.4	97.6	100.0	4.2	19.6	85.6		4.8	9.2	37.4
	0.2	7.2	97.4	100.0	4.0	18.8	84.6		4.8	9.2	35.2

surprising, since the "optimal" sampling in SAS takes account of the data structure like our theoretical analysis in Section 2.3. The US usually owns a higher power compared to US₀ due to the test statistic of the former one based on a larger set with sample size $n + n_0$. However, US performs a little poorly in terms of the empirical size than US₀, since the samples for the tests are dependent on the estimated θ in the pilot study.

In Figure 3, we further compare the proposed method SAS with Guo from Guo et al. (2016) and SZ from Stute and Zhu (2002) under Scenario I. In most cases, our method SAS performs effectively and leads to a higher power than the competitors. The

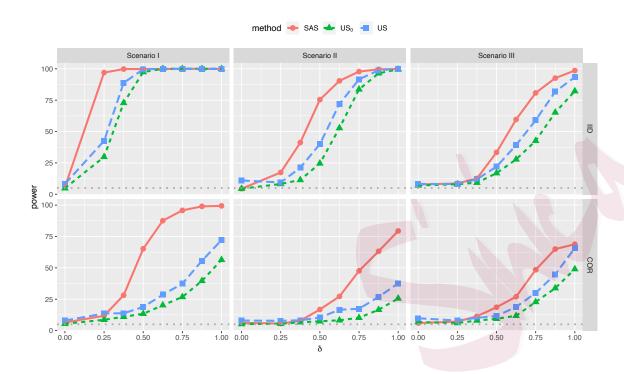


Figure 2: Empirical sizes and powers (%) for SAS, US₀ and US under Scenarios I–III.

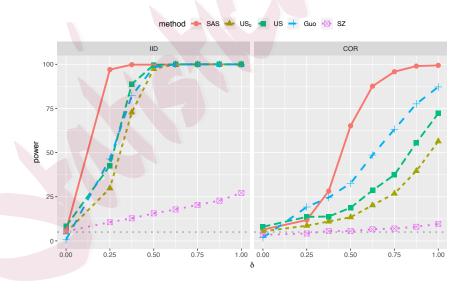


Figure 3: Empirical sizes and powers (%) for SAS, US₀, US, Guo in Guo et al. (2016) and SZ in Stute and Zhu (2002) under Scenario I.

SZ performs not well. This is not surprising, since the local-smoothing-based methods work better than the global testing procedure SZ for local alternatives as suggested by existing numerical studies in the literature. We also observe that the SAS test tends to be more conservative than the Guo method slightly when the signal δ is very small, especially under the COR case. This can be understood that the Guo uses the $n + n_0$ samples for testing. However, the power of SAS increases quickly as δ increases and its adaptive-sampling advantage becomes remarkable.

In what follows, we consider one general varying coefficient model.

• Scenario IV (Varying coefficient model): $G(\mathbf{X}; \boldsymbol{\beta}, \mathbf{g}) = \boldsymbol{\beta}_1(X_1)X_2 + \boldsymbol{\beta}_2(X_1)X_3 + \delta \cdot 0.3 \left(\boldsymbol{\theta}^{\top}\mathbf{X}\right)^3$ under p = 4, $\boldsymbol{\beta}_1(X_1) = \sin(X_1) + \cos(X_1)$, $\boldsymbol{\beta}_2(X_1) = 2X_1(1 - X_1)$, and $\boldsymbol{\theta} = (0, 1, 1, 1)^{\top} / \sqrt{3}$. $X_1 \sim \mathcal{U}(0, 1)$, and $X_2, X_3, X_4 \sim \mathcal{N}(0, \mathbf{I}_p)$, where $\mathcal{U}(0, 1)$ is the standard uniform distribution.

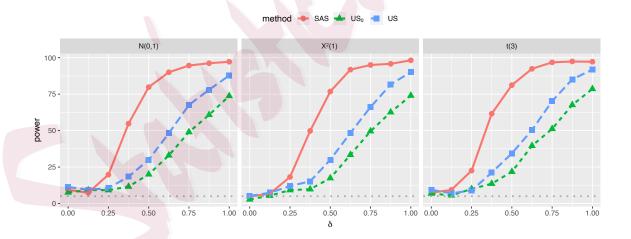


Figure 4: Empirical sizes and powers (%) for SAS, US₀ and US under Scenario IV.

In Scenario IV, we examine the robustness of the proposed test when the errors ε_i 's are from different distributions. Figure 4 shows that the empirical sizes and powers

for SAS, US₀ and US under three standardized error distributions, i.e., $\mathcal{N}(0,1)$, $\chi^2(1)$ and t(3). Again, it can be seen that the SAS outperforms the US₀ and US uniformly and there are no significant differences among the power curves under different error distributions.

Table 2: Computing times (seconds) for SAS, US₀ and US when ε_i 's are i.i.d from $\mathcal{N}(0,1)$ and the full sample size N is from 10^3 to 10^6 under Scenario IV.

δ	N Method	10^{3}	10^4	10^{5}	10^6
	SAS	0.237	0.787	3.913	29.907
0.00	US_0	0.215	0.251	0.239	0.259
	US	0.414	0.471	0.310	0.350
	FULL	0.426	7.794	750.735	77033.543
	SAS	0.233	0.487	3.405	30.336
0.25	US_0	0.196	0.222	0.230	0.252
	US	0.268	0.316	0.307	0.330
	FULL	0.432	7.845	790.955	78946.865
	SAS	0.233	0.484	3.407	30.164
0.5	US_0	0.196	0.222	0.234	0.254
	US	0.267	0.316	0.305	0.324
	FULL	0.433	7.832	795.734	81033.674

Finally, to further investigate the computational benefit of our SAS procedure in large-scale datasets, the computing time is reported in Table $\boxed{2}$. As the full sample size N increases, subsampling methods take significantly less computing time compared to the full data approach. We also observe that the computing time of the SAS procedure roughly linearly increases as the sample size N increases but the time based on the full sample is quadratically increased, which is consistent with our theoretical computing cost analysis. It is not surprising that the US₀ and US run fast since there is no sampling distribution need to be estimated.

4.2 Real data analysis

Wave Energy Converters (WECs) are of interest to governments and industry as a way of complementing other renewable energy sources (Neshat et al., 2018) since they have advantages in terms of high availability of resource. However, huge information for the WECs can be recorded and monitored with buoys, but analyzing all the collected data will lead to a heavy computing burden. To overcome this challenge, we apply the proposed SAS to the "Wave Energy Converters Data Set". This data includes all 216000 measuring sample points under the three real wave scenarios from the southern coast of Australia and it is available from UCI Machine Learning Repository http://archive.ics.uci.edu/ml/datasets/Wave+Energy+Converters.

The interest in this problem is to study the relationships between the total absorbed wave power output (Y) and the WECs positions (\mathbf{X}) . Here we focus on the first three WECs' positions, i.e. p=6. Since the power output is often positive but large, we take the logarithm transformation to the response Y. We randomly select 70% of the

full data as the training set (N=151200) and let the rest as the validation set. For the pilot study, $n_0=300$ samples are uniformly sampled from the training set and are used to estimate the projection direction $\boldsymbol{\theta}$ and optimal sampling distribution $f(\cdot)$.

In this application, we would like to check whether the linear model (LM) and mult-index model (MIM) are sufficient to describe the relationship between response and covariates. Specifically, we denote the multi-index model with two indexes as MIM2. The MIM3 and MIM4 are defined similarly. In Table $\boxed{3}$, we compare the estimated p-values of SAS, US₀ and US under different model assumptions.

Table 3: The estimated p-value and MSPE of wave energy converters dataset.

		MSPE		
Model	SAS	\mathbf{US}_0	US	
LM	0.000	0.000	0.000	164.686
MIM2	0.005	0.097	0.055	0.201
$_{ m MIM3}$	0.970	0.444	0.848	0.195
MIM4	0.928	0.607	0.525	0.200

It is clear that all three methods suggest a nonlinear relationship between wave power and the WECs positions due to p-values equal to 0 under the linear model assumption. Consider the significant level α is 0.05. We find that US₀ and US perform similarly and both of them tend to choose MIM2 for model construction. However, SAS suggests MIM3 is more reliable, as it owns a much smaller p-value 0.005 than the other two methods under the model assumption MIM2. This is probably because the full data consists of three real wave scenarios. The result can be further verified by the

comparison of the mean square prediction errors (MSPE) on the test data. In the last column of Table 3, the MSPEs under MIM3 are much smaller than the other models.

5. Concluding remarks

Model checking in large-scale datasets is an important preliminary step for statistical analysis and machine learning. In this paper, we construct a *structure-adaptive-sampling*, SAS test, in a general semiparametric framework to overcome the large-scale dataset computational bottleneck with limited budget or resources. The SAS procedure could select those most informative samples with an optimal sampling procedure. It is shown that our proposed method can asymptotically achieve locally best power. The asymptotic and numerical results demonstrate the advantages of the proposed procedure in testing and computation via comparing to the uniform sampling strategy and some existing model checking approaches.

In general, the SAS procedure can be readily extended to many other testing problems with sampling techniques in modern large-scale datasets analysis, such as clustering several regression curves or datasets, as long as the design point sampling is allowed in the process. Our analysis shows that the covariate correlation plays an important role in the performance of the SAS procedure. Though our results reveal that the superiority of the proposed method is valid under certain correlations, it is of interest to see how to incorporate the correlation information into the testing procedure in an efficient way. Besides, we only use a small pilot dataset to estimate the projection direction θ when no more information is given. It requires more research to obtain a

more accurate estimation of the projection direction for test power improvement.

Acknowledgements

The authors thank the editor, the associate editor, and the anonymous referees for their valuable comments that improved the paper significantly. The authors also thank Professor Changliang Zou for his constructive suggestions and unselfish assistance. Han and Wang's research were supported by the National Natural Science Foundation of China Grants No.11931001, 11690015, 11925106, 11971247, the National Science Foundation of Tianjin 18JCJQJC46000, and the Fundamental Research Funds for the Central Universities of Nankai University 63201162. Ma's research was partially supported by the U.S. National Science Foundation under grants DMS-1903226, DMS-1925066, and the U.S. National Institute of Health under grant R01GM122080. Ren's work was partially sponsored by Shanghai Sailing Program.

A. Appendix: Proofs

In this Appendix, we prove the main theoretical results in our paper. Before we present the proofs of the main results, we first state several essential lemmas whose proofs can be founded in the Supplementary Material. Denote $\Upsilon_j^* = G(\mathbf{X}_j; \widehat{\boldsymbol{\beta}}, \widehat{\mathbf{g}}) - G(\mathbf{X}; \boldsymbol{\beta}^*, \mathbf{g}^*)$ and $\boldsymbol{\Upsilon}^* = (\Upsilon_1^*, \dots, \Upsilon_n^*)^{\top}$. Under \mathbb{H}_0 , $\Upsilon_j^* = G(\mathbf{X}_j; \widehat{\boldsymbol{\beta}}, \widehat{\mathbf{g}}) - G(\mathbf{X}; \boldsymbol{\beta}_0, \mathbf{g}_0)$ due to $(\boldsymbol{\beta}^*, \mathbf{g}^*) = (\boldsymbol{\beta}_0, \mathbf{g}_0)$. For the notation simplicity, we denote the following form as

$$W_n(\mathbf{s}, \mathbf{t}) = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n \frac{s_i t_j}{\sqrt{f(\omega_i) f(\omega_j)}} K_h(\omega_i - \omega_j),$$

where $\mathbf{s} = (s_1, \dots, s_n)^{\top}$ and $\mathbf{t} = (t_1, \dots, t_n)^{\top}$ are two sequences. For example, our test statistic can be written as $V_f(\boldsymbol{\theta}) = W_n(\widehat{\boldsymbol{\varepsilon}}, \widehat{\boldsymbol{\varepsilon}})$, where $\widehat{\boldsymbol{\varepsilon}} = (\widehat{\varepsilon}_1, \dots, \widehat{\varepsilon}_n)^{\top}$, $\widehat{\varepsilon}_i = Y_i - G(\mathbf{X}_i; \widehat{\boldsymbol{\beta}}, \widehat{\mathbf{g}})$.

Lemma A.1. Suppose the Assumptions 1—3 and 6 hold, and denote $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^{\top}$.

With a given $\boldsymbol{\theta}$, then under \mathbb{H}_0 we have $nh^{1/2}W_n(\boldsymbol{\varepsilon}, \boldsymbol{\varepsilon})/\sigma_V \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$.

Lemma A.2. Suppose the Assumptions 4 and 5 hold, then

$$\sup_{\nu \in \Gamma} \left| G\left(\mathbf{X}; \widehat{\boldsymbol{\beta}}, \widehat{\mathbf{g}}(\nu) \right) - G\left(\mathbf{X}; \boldsymbol{\beta}^*, \mathbf{g}^*(\nu) \right) \right| = O_p \left\{ b^2 + n^{-1/2} + (nb/\log n)^{-1/2} \right\}.$$

Lemma A.3. Suppose the Assumptions hold. Given $\boldsymbol{\theta}$, then under \mathbb{H}_0 we have: (i) $W_n(\boldsymbol{\varepsilon}, \boldsymbol{\Upsilon}^*) = o_p(n^{-1}h^{-1/2})$; (ii) $W_n(\boldsymbol{\Upsilon}^*, \boldsymbol{\Upsilon}^*) = o_p(n^{-1}h^{-1/2})$.

Lemma A.4. Suppose the Assumptions $\boxed{1}$ — $\boxed{6}$ hold. Given $L(\cdot)$ is a continuously differentiable function, which satisfies $|L(\mathbf{X})| \leq \varphi(\mathbf{X})$ for all $\mathbf{X} \in \mathbb{R}^p$ and $\mathbb{E} \{\varphi^2(\mathbf{X})\} < \infty$. Denote $\mathbf{L} = \{L(\mathbf{X}_i)\}_{i=1}^n$. With a given $\boldsymbol{\theta}$, then under \mathbb{H}_0 the following result holds

$$W_n(\boldsymbol{\varepsilon}, \mathbf{L}) = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n \frac{\varepsilon_i L(\mathbf{X}_j)}{\sqrt{f(\omega_i)f(\omega_j)}} K_h(\omega_i - \omega_j) = O_p(n^{-1/2}).$$

Lemma A.5. Suppose Assumptions 1–3 hold. With a given $\boldsymbol{\theta}$, then for the density estimator $\widehat{f}(\omega)$ in (2.5) from the pilot study, we have

$$\sup_{\omega \in \Omega} |\widehat{f}(\omega) - f(\omega)| = O_p \left(h_f^2 + \sqrt{\frac{\log n_0}{n_0 h_f}} \right).$$

We give a sketch of our proofs. We first derive the asymptotic distribution of our SAS test statistic with a given dimension reduction direction $\boldsymbol{\theta}$. Next, the method of deriving the optimal sampling distribution $f(\cdot)$ is conducted in Theorem 1. Then we

give the power function for our SAS test in Theorem 2. At last, we extend our sampling strategy to the "singular" signal case in Corollary 1. The relevant proof of estimated dimension reduction direction is delineated in the Supplementary Material.

Proposition A.1. Suppose Assumptions [1]—[6] hold. Under the local alternatives \mathbb{H}_{1n} (2.3) with $\delta_n = (nh^{1/2})^{-1/2}$ and a given $\boldsymbol{\theta}$, we have $T_f - \mathbb{E}_f \{l^2(\omega)\} / \sigma_V \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$.

Proof. our SAS test statistic $T_f = nh^{1/2}V_f(\boldsymbol{\theta})$ and $V_f(\boldsymbol{\theta})$ has the following decomposition

$$V_{f}(\boldsymbol{\theta}) = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i}^{n} \frac{\epsilon_{i}}{\sqrt{f(\omega_{i})}} K_{h}(\omega_{i} - \omega_{j}) \frac{\epsilon_{j}}{\sqrt{f(\omega_{j})}}$$

$$- \frac{2}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i}^{n} \frac{\epsilon_{i}}{\sqrt{f(\omega_{i})}} K_{h}(\omega_{i} - \omega_{j}) \frac{\Upsilon_{j}^{*}}{\sqrt{f(\omega_{j})}}$$

$$+ \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i}^{n} \frac{\Upsilon_{i}^{*}}{\sqrt{f(\omega_{i})}} K_{h}(\omega_{i} - \omega_{j}) \frac{\Upsilon_{j}^{*}}{\sqrt{f(\omega_{j})}}$$

$$=: W_{n}(\boldsymbol{\epsilon}, \boldsymbol{\epsilon}) - 2W_{n}(\boldsymbol{\epsilon}, \boldsymbol{\Upsilon}^{*}) + W_{n}(\boldsymbol{\Upsilon}^{*}, \boldsymbol{\Upsilon}^{*}),$$

where $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^{\top}$ and $\epsilon_i = Y_i - G(\mathbf{X}_i; \boldsymbol{\beta}^*, \mathbf{g}^*) = \delta_n l(\omega_i) + \varepsilon_i$ under \mathbb{H}_{1n} .

Similarly, we can write the first term as $W_n(\boldsymbol{\epsilon}, \boldsymbol{\epsilon}) = W_n(\boldsymbol{\epsilon}, \boldsymbol{\epsilon}) + 2\delta_n W_n(\boldsymbol{\epsilon}, \boldsymbol{l}) + \delta_n^2 W_n(\boldsymbol{l}, \boldsymbol{l})$ where $\boldsymbol{l} = (l(\omega_1), \dots, l(\omega_n))^{\top}$. Note that $W_n(\boldsymbol{l}, \boldsymbol{l})$ is a U-statistic of order two with kernel $H_n(\omega_i, \omega_j)$, say $H_n(\omega_i, \omega_j) = l(\omega_i)l(\omega_j)K_h(\omega_i - \omega_j)\left\{f(\omega_i)f(\omega_j)\right\}^{-1/2}$, and

$$\mathbb{E}\left\{H_n(\omega_1, \omega_2)\right\} = \mathbb{E}\left[\mathbb{E}\left\{H_n(\omega_1, \omega_2)|\omega_1\right\}\right]$$

$$= \frac{1}{h} \int K(u)l(\omega)l(\omega - hu)\sqrt{f(\omega)f(\omega - hu)}h du d\omega$$

$$= \int K(u)l^2(\omega)f(\omega)du d\omega + o(1)$$

$$= \mathbb{E}_f\left\{l^2(\omega)\right\} + o(1).$$

$$\mathbb{E}\left\{H_n^2(\omega_1, \omega_2)\right\} = \frac{1}{h^2} \int \frac{l^2(\omega_i)l^2(\omega_j)}{f(\omega_i)f(\omega_j)} K^2 \left(\frac{\omega_i - \omega_j}{h}\right) f(\omega_i) f(\omega_j) d\omega_i d\omega_j$$

$$= \frac{1}{h} \int K^2(u) du \cdot \int l^2(\omega_j + hu) l^2(\omega_j) d\omega_j$$

$$= \frac{1}{h} \int K^2(u) du \cdot \int l^4(\omega_j) d\omega_j + o(1/h)$$

$$= O(1/h) = o(n).$$

By Lemma S.2, we can get $W_n(\boldsymbol{l},\boldsymbol{l}) = \mathbb{E}_f \{l^2(\omega)\} + o_p(1)$. Combining this with the results in Lemma A.1 and Lemma A.3, we have $nh^{1/2}W_n(\boldsymbol{\epsilon},\boldsymbol{\epsilon})/\sigma_V - \mathbb{E}_f \{l^2(\omega)\}/\sigma_V \xrightarrow{\mathcal{L}} \mathcal{N}(0,1)$.

Next, we consider the second term $W_n(\epsilon, \Upsilon^*) = W_n(\epsilon, \Upsilon^*) + \delta_n W_n(l, \Upsilon^*)$. By Lemma A.2 and

$$\frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i}^{n} \frac{l(\omega_i)}{\sqrt{f(\omega_i)}} \frac{1}{\sqrt{f(\omega_j)}} K_h(\omega_i - \omega_j) = \mathbb{E}_f \left\{ l(\omega) \right\} + o_p(1),$$

we can claim that $W_n(\mathbf{l}, \Upsilon^*) = O_p(n^{-1/2})O_p\{b^2 + n^{-1/2} + (nb/\log n)^{-1/2}\}$. As a consequence, we have $W_n(\epsilon, \Upsilon^*) = o_p(n^{-1}h^{-1/2})$ by Lemma A.3.

Finally, it can be checked that $W_n(\Upsilon^*, \Upsilon^*) = o_p(n^{-1}h^{-1/2})$ by Lemma A.3, from which the assertion of the proposition holds.

Proof of Theorem 1

By Proposition A.1, it implies that the asymptotic power function of proposed test based on T_f only depends on $\mathbb{E}_f \{l^2(\omega)\}/\sigma_V$. Note that $\sigma_V^2 = 2\sigma^4 |\Omega| \int K^2(u) du$ is a constant not depending on the density $f(\omega)$ under a prescribed kernel function $K(\cdot)$. According to Cauchy-Schwarz inequality, we have

$$\mathbb{E}_f\left\{l^2(\omega)\right\} = \int l^2(\omega)f(\omega)d\omega \le \left\{\int l^4(\omega)d\omega\right\}^{1/2} \left\{\int f^2(\omega)d\omega\right\}^{1/2},$$

and the equality holds if and only if $f(\omega) = \rho l^2(\omega)$, where ρ is some constant. Thus, if the sampling distribution $f(\omega)$ is proportional to $l^2(\omega)$, i.e. $f(\omega) = l^2(\omega)/\int l^2(\omega)d\omega$, we can maximize the asymptotic power function. That is the locally best power will be reached by choosing $f(\omega) = l^2(\omega)/\int l^2(\omega)d\omega$.

Proof of Theorem 2

The test statistic is $T_{\widehat{f}}$ with estimated $\widehat{f}(\omega)$ from the pilot study. Under the local alternatie \mathbb{H}_{1n} , the power is $\Pi_f = \Pr(T_{\widehat{f}} > z_{\alpha})$. By Proposition A.1, we have

$$\Pr(T_{\widehat{f}} > z_{\alpha}) - \Phi\left(-z_{\alpha} + \mathbb{E}_{\widehat{f}}\{l^{2}(\omega)\}/\sigma_{V}\right) \xrightarrow{P} 0.$$

Using the uniform convergence rate given in Lemma A.5, we have

$$\mathbb{E}_{\widehat{f}}\left\{l^2(\omega)\right\} - \int l^4(\omega) d\omega / \int l^2(\omega) d\omega \stackrel{P}{\to} 0,$$

provided that the order of signal of strength is larger than that the maximum noise level say $(nh^{1/2})^{-1/2}/\sqrt{\log n_0/n_0h_f} \to \infty$. The condition $(n_0/n)(h_f^2/h)^{1/2}/(\log n_0)^c \to \infty$ implies that the result holds.

Proof of Corollary 1.

For simplicity, we assume there is only one signal region Ω_n , and the proof for the case with more than one region is similar. By condition $(n_0/n)(h_f^2/h)^{1/2}/\left\{a_n^{1/2}(\log n_0)^c\right\} \to \infty$, it suffices to show that the Nadaraya-Watson estimator of $M(\omega)$ in (2.4) is still a uniformly consistent one except for the boundary under "singular" local alternatives (2.6) (Ren et al., 2020), say

$$\sup_{\omega \in \Omega_n \setminus \mathbb{B}_h} \left| \widehat{M}(\omega) - M(\omega) \right| = O_p \left(h_f^2 \delta_n' + \sqrt{\frac{a_n \log n_0}{n_0 h_f}} \right), \sup_{\omega \in (\Omega \setminus \Omega_n) \setminus \mathbb{B}_h} \left| \widehat{M}(\omega) \right| = O_p \left(\sqrt{\frac{a_n \log n_0}{n_0 h_f}} \right),$$

SAS MODEL CHECKING

where \mathbb{B}_h denotes a one-dimensional interval with radius h that is around the boundary of Ω_n . The rate of bias term is obvious as $\epsilon_i = \delta_n l(\omega_i) + \varepsilon_i$ under (2.6), and the rate of variance term is reasonable because it uniformly takes maximum over $\omega \in \Omega_n$. Then we have

$$\int_{\Omega_n \cup \mathbb{B}_h} f(\omega) d\omega = 1 + O_p\left(\sqrt{\frac{a_n \log n_0}{n_0 h_f}}\right), \quad \Pr\left\{f(\omega_i) < \xi_{n_0}^{1/2}, \forall \omega_i \notin \Omega_n \cup \mathbb{B}_h\right\} \to 1,$$

where ξ_{n_0} is the threshold defined in (2.5). Hence, we have

$$\Pr\left(T_{\widehat{f}} > z_{\alpha} | \mathcal{X}\right) - \Phi\left(\mu(f, \Omega_n) / \sigma_V(f, \Omega_n) - z_{\alpha}\right) \xrightarrow{P} 0,$$

where $\sigma_V^2(f,\Omega_n) \approx 2\sigma^4 |\Omega_n| \int K^2(u) du$, and

$$\mu(f, \Omega_n) = nh^{1/2} {\delta'_n}^2 \int_{\Omega_n \cup \mathbb{B}_h} l^2(\omega) f(\omega) d\omega$$

$$\approx nh^{1/2} {\delta'_n}^2 \int_{\Omega_n \setminus \mathbb{B}_h} l^4(\omega) d\omega / \int_{\Omega_n \setminus \mathbb{B}_h} l^2(\omega) d\omega$$

$$\geq nh^{1/2} {\delta'_n}^2 \int_{\Omega_n} l^2(\omega) l_{\min}^2 d\omega / \int_{\Omega_n} l^2(\omega) d\omega$$

$$= a_n^{-1/2} l_{\min}^2.$$

Thus, the power of our SAS is not small than $\Phi(\mu(f,\Omega_n)/\sigma_V(f,\Omega_n)-z_\alpha)$. Finally, by the assumption that $a_n \to 0$, the asymptotic power converges to 1.

B. Supplementary Material

- A Supplementary PDF contains the proofs of several technical lemmas, the relevant t proof of estimated dimension reduction direction, and some additional simulation results.
- R computing code for the numerical analysis.

References

- Ai, M., Yu, J., Zhang, H., and Wang, H. (2021). Optimal subsampling algorithms for big data regressions. *Statistica Sinica*, **6**:363–392.
- Balakrishnan, S. and Madigan, D. (2008). Algorithms for sparse linear classifiers in the massive data setting. *Journal of Machine Learning Research*, **9**:313–337.
- Battey, H., Fan, J., Liu, H., Lu, J., and Zhu, Z. (2018). Distributed testing and estimation under sparse high dimensional models. *The Annals of Statistics*, **46**(3):1352–1382.
- Cook, R. D. and Forzani, L. (2009). Likelihood-based sufficient dimension reduction. Journal of the American Statistical Association, 104(485):197–208.
- Cook, R. D. and Weisberg, S. (1991). Sliced inverse regression for dimension reduction:

 Comment. Journal of the American Statistical Association, 86(414):328–332.
- Dette, H. (1999). A consistent test for the functional form of a regression based on a difference of variance estimators. *The Annals of Statistics*, **27**(3):1012–1040.
- Fan, J. and Huang, L.-S. (2001). Goodness-of-fit tests for parametric regression models.

 Journal of the American Statistical Association, 96(454):640–652.
- Fan, J. and Huang, T. (2005). Profile likelihood inferences on semiparametric varying-coefficient partially linear models. *Bernoulli*, **11**(6):1031–1057.
- Fan, J., Zhang, C., and Zhang, J. (2001). Generalized likelihood ratio statistics and wilks phenomenon. *The Annals of Statistics*, **29**(1):153–193.

- Fan, J. and Zhang, W. (1999). Statistical estimation in varying coefficient models. The Annals of Statistics, 27(5):1491–1518.
- Ferragut, E. M. and Laska, J. (2012). Randomized sampling for large data applications of svm. In 2012 11th International Conference on Machine Learning and Applications, volume 1, pages 350–355.
- González-Manteiga, W. and Crujeiras, R. M. (2013). An updated review of goodness-of-fit tests for regression models. *Test*, **22**(3):361–411.
- Guerre, E. and Lavergne, P. (2005). Data-driven rate-optimal specification testing in regression models. *The Annals of Statistics*, **33**(2):840–870.
- Guo, X., Wang, T., and Zhu, L. (2016). Model checking for parametric single-index models: a dimension reduction model-adaptive approach. *Journal of the Royal Statistical Society:*Series B (Statistical Methodology), 78(5):1013–1035.
- Guo, X. and Zhu, L. (2017). A review on dimension-reduction based tests for regressions. In From Statistics to Mathematical Finance, pages 105–125. Springer.
- Hardle, W., Hall, P., and Ichimura, H. (1993). Optimal smoothing in single-index models.

 The Annals of Statistics, 21(1):157–178.
- Hardle, W. and Mammen, E. (1993). Comparing nonparametric versus parametric regression fits. *The Annals of Statistics*, **21**(4):1926–1947.
- Hart, J. D. (1997). Nonparametric Smoothing and Lack-of-Fit Tests. Springer Series in Statistics. Springer New York, New York, NY.

- Hastie, T. and Tibshirani, R. (1986). Generalized additive models. *Statistical Science*, **1**:297–318.
- Hastie, T. and Tibshirani, R. (1993). Varying-coefficient models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **55**(4):757–779.
- Horowitz, J. L. and Mammen, E. (2004). Nonparametric estimation of an additive model with a link function. *The Annals of Statistics*, **32**(6):2412–2443.
- Ichimura, H. (1993). Semiparametric least squares (sls) and weighted sls estimation of single-index models. *Journal of Econometrics*, **58**(1-2):71–120.
- Jordan, M. I., Lee, J. D., and Yang, Y. (2019). Communication-efficient distributed statistical inference. *Journal of the American Statistical Association*, **114**(526):668–681.
- Kleiner, A., Talwalkar, A., Sarkar, P., and Jordan, M. I. (2014). A scalable bootstrap for massive data. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 76(4):795–816.
- Li, B. and Wang, S. (2007). On directional regression for dimension reduction. *Journal of the American Statistical Association*, **102**(479):997–1008.
- Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, **86**(414):316–327.
- Ma, P., Mahoney, M. W., and Yu, B. (2015). A statistical perspective on algorithmic leveraging. *Journal of Machine Learning Research*, **16**(1):861–911.

- Ma, P. and Sun, X. (2015). Leveraging for big data regression. Wiley Interdisciplinary Reviews: Computational Statistics, 7(1):70–76.
- Neshat, M., Alexander, B., Wagner, M., and Xia, Y. (2018). A detailed comparison of metaheuristic methods for optimising wave energy converter placements. In *Proceedings of the* Genetic and Evolutionary Computation Conference, pages 1318–1325.
- Ren, H., Zou, C., Chen, N., and Li, R. (2020). Large-scale datastreams surveillance via pattern-oriented-sampling. *Journal of the American Statistical Association*, pages 1–15.
- Schifano, E. D., Wu, J., Wang, C., Yan, J., and Chen, M.-H. (2016). Online updating of statistical inference in the big data setting. *Technometrics*, **58**(3):393–403.
- Speckman, P. (1988). Kernel smoothing in partial linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **50**(3):413–436.
- Stute, W., Manteiga, W. G., and Quindimil, M. P. (1998). Bootstrap approximations in model checks for regression. *Journal of the American Statistical Association*, **93**(441):141–149.
- Stute, W. and Zhu, L.-X. (2002). Model checks for generalized linear models. *Scandinavian Journal of Statistics*, **29**(3):535–545.
- Stute, W., Zhu, L.-X., et al. (2005). Nonparametric checks for single-index models. *The Annals of Statistics*, **33**(3):1048–1083.
- Tan, F., Zhu, X., and Zhu, L. (2018). A projection-based adaptive-to-model test for regressions. Statistica Sinica, pages 157–188.

- Tur, G., Hakkani-Tür, D., and Schapire, R. E. (2005). Combining active and semi-supervised learning for spoken language understanding. *Speech Communication*, **45**(2):171–186.
- Wang, H., Yang, M., and Stufken, J. (2019). Information-based optimal subdata selection for big data linear regression. *Journal of the American Statistical Association*, **114**(525):393–405.
- Wang, H., Zhu, R., and Ma, P. (2018). Optimal subsampling for large sample logistic regression. *Journal of the American Statistical Association*, **113**(522):829–844.
- Wang, Y., Yu, A. W., and Singh, A. (2017). On computationally tractable selection of experiments in measurement-constrained regression models. *Journal of Machine Learning Research*, **18**(1):5238–5278.
- Xia, Y. (2009). Model checking in regression via dimension reduction. *Biometrika*, **96**(1):133–148.
- Xia, Y., Tong, H., Li, W., and Zhu, L.-X. (2002). An adaptive estimation of dimension reduction space. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 64(3):363–410.
- Yao, Y. and Wang, H. (2019). Optimal subsampling for softmax regression. *Statistical Papers*, **60**(4):235–249.
- Yu, J., Wang, H., Ai, M., and Zhang, H. (2020). Optimal distributed subsampling for maximum quasi-likelihood estimators with massive data. *Journal of the American Statistical Association*, pages 1–12.

- Zhao, Y., Zou, C., and Wang, Z. (2019). A scalable nonparametric specification testing for massive data. *Journal of Statistical Planning and Inference*, **200**:161–175.
- Zheng, J. X. (1996). A consistent test of functional form via nonparametric estimation techniques. *Journal of Econometrics*, **75**(2):263–289.
- Zhu, X., Guo, X., and Zhu, L. (2017). An adaptive-to-model test for partially parametric single-index models. *Statistics and Computing*, **27**(5):1193–1204.