

This is a repository copy of *Semiparametric Ultra-High Dimensional Model Averaging of Nonlinear Dynamic Time Series*.

White Rose Research Online URL for this paper:  
<http://eprints.whiterose.ac.uk/114572/>

Version: Accepted Version

---

**Article:**

Chen, Jia [orcid.org/0000-0002-2791-2486](https://orcid.org/0000-0002-2791-2486), Li, Degui [orcid.org/0000-0001-6802-308X](https://orcid.org/0000-0001-6802-308X), Linton, Oliver et al. (1 more author) (2018) Semiparametric Ultra-High Dimensional Model Averaging of Nonlinear Dynamic Time Series. *Journal of the American Statistical Association*. pp. 919-932. ISSN 0162-1459

<https://doi.org/10.1080/01621459.2017.1302339>

---

**Reuse**

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Semiparametric Ultra-High Dimensional Model Averaging of Nonlinear Dynamic Time Series

Jia Chen\*      Degui Li†      Oliver Linton‡      Zudi Lu§

Version: February 22, 2017

## Abstract

We propose two semiparametric model averaging schemes for nonlinear dynamic time series regression models with a very large number of covariates including exogenous regressors and autoregressive lags. Our objective is to obtain more accurate estimates and forecasts of time series by using a large number of conditioning variables in a nonparametric way. In the first scheme, we introduce a Kernel Sure Independence Screening (KSIS) technique to screen out the regressors whose marginal regression (or auto-regression) functions do not make a significant contribution to estimating the joint multivariate regression function; we then propose a semiparametric penalized method of Model Averaging MArginal Regression (MAMAR) for the regressors and auto-regressors that survive the screening procedure, to further select the regressors that have significant effects on estimating the multivariate regression function and predicting the future values of the response variable. In the second scheme, we impose an approximate factor modelling structure on the ultra-high dimensional exogenous regressors and use the principal component analysis to estimate the latent common factors; we then apply the penalized MAMAR method to select the estimated common factors and the lags of the response variable that are significant. In each of the two schemes, we construct the optimal combination of the significant marginal regression and auto-regression functions. Asymptotic properties for these two schemes are derived under some regularity conditions. Numerical studies including both simulation and an empirical application to forecasting inflation are given to illustrate the proposed methodology.

*Keywords:* Kernel smoother, penalized MAMAR, principal component analysis, semiparametric approximation, sure independence screening, ultra-high dimensional time series.

---

\*Department of Economics and Related Studies, University of York, YO10 5DD, UK. E-mail: [jia.chen@york.ac.uk](mailto:jia.chen@york.ac.uk)

†Department of Mathematics, University of York, YO10 5DD, UK. E-mail: [degui.li@york.ac.uk](mailto:degui.li@york.ac.uk).

‡Faculty of Economics, University of Cambridge, Austin Robinson Building, Sidgwick Avenue, CB3 9DD, UK. E-mail: [ob120@cam.ac.uk](mailto:ob120@cam.ac.uk).

§Statistical Sciences Research Institute and School of Mathematical Sciences, University of Southampton, SO17 1BJ, UK. E-mail: [Z.Lu@soton.ac.uk](mailto:Z.Lu@soton.ac.uk).

# 1 Introduction

Nonlinear time series modelling taking account of both dynamic lags of response variable and exogenous regressors is of wide interest in applications. We suppose that  $Y_t$ ,  $t = 1, \dots, n$ , are  $n$  observations collected from a stationary time series process, and often we are interested in the multivariate dynamic regression function

$$m(\mathbf{x}) = \mathbf{E}(Y_t | \mathbf{X}_t = \mathbf{x}), \quad (1.1)$$

where  $Y_t$  is the response variable, and  $\mathbf{X}_t = (\mathbf{Z}_t^\top, \mathbf{Y}_{t-1}^\top)^\top$  with  $\mathbf{Z}_t = (Z_{t1}, Z_{t2}, \dots, Z_{tp_n})^\top$  and  $\mathbf{Y}_{t-1} = (Y_{t-1}, Y_{t-2}, \dots, Y_{t-d_n})^\top$  being a  $p_n$ -dimensional vector consisting of exogenous regressors and a vector of  $d_n$  lags of  $Y_t$ , respectively. Here the superscript  $\top$  stands for the transpose of a vector (or a matrix). We allow that both  $p_n$  and  $d_n$  could increase with the sample size  $n$ , and that  $\mathbf{Z}_t$  could include lags of the exogenous regressors and has a large dimension  $p_n$ , allowed to be even larger than the sample size  $n$ . Such an ultra-high dimensional time series setting poses a great challenge in estimating the regression function  $m(\mathbf{x})$  and the subsequent forecast of the response.

When the dimension of  $\mathbf{X}_t$  is low (say 1 or 2), it is well known that the conditional regression function  $m(\mathbf{x})$  can be well estimated by using some commonly-used nonparametric methods such as the kernel method, the local polynomial method, and the spline method (c.f., [Green and Silverman, 1994](#); [Wand and Jones, 1995](#); [Fan and Gijbels, 1996](#)). However, if  $\mathbf{X}_t$  is of large dimension, owing to the so-called ‘‘curse of dimensionality’’, a direct use of nonparametric methods leads to a very poor estimation and forecasting result. Hence, various nonparametric and semiparametric models, such as additive models, varying coefficient models and partially linear models, have been proposed to deal with the curse of dimensionality (c.f., [Teräsvirta et al, 2010](#)). A recent paper by [Li et al \(2015\)](#) develops a flexible semiparametric forecasting model, termed ‘‘*Model Averaging MArginal Regression*’’ (MAMAR). It seeks to optimally combine nonparametric low-dimensional marginal regressions, which helps to improve the accuracy of predicting future values of the nonlinear time series.

The idea of the model averaging approach is to combine several candidate models by assigning higher weights to better candidate models. Under the linear regression setting with the dimension of covariates smaller than the sample size, there has been an extensive literature on various model averaging methods, see, for example, the AIC and BIC model averaging ([Akaike, 1979](#); [Raftery et al, 1997](#); [Claeskens and Hjort, 2008](#)), the Mallows  $C_p$  model averaging ([Hansen, 2007](#); [Wan et al, 2010](#)), and the jackknife model averaging ([Hansen and Racine, 2012](#)). However, in the case of ultra-high dimensional time series, these methods may not perform well and the associated asymptotic theory fails. To address this issue, [Ando and Li \(2014\)](#) propose a two-step model averaging method for a high-dimensional linear regression with the dimension of the covariates larger than the sample size and show that such a method works well both theoretically and numerically. Recently [Cheng and Hansen \(2015\)](#) study the model averaging of the factor-augmented linear regression by applying the principal component analysis on the high-dimensional covariates to estimate the unobservable factor regressors.

In this paper, our main objective is to propose semiparametric ultra-high dimensional model averaging schemes for studying the nonlinear dynamic regression structure for (1.1), which generalizes the existing approaches. On the one hand, we relax the restriction of linear modelling assumed in Ando and Li (2014) and Cheng and Hansen (2015), and on the other hand, we extend the recent work of Li *et al* (2015) to the ultra high dimensional case, thereby providing a much more flexible framework for nonlinear dynamic time series forecasting.

Throughout the paper, we assume that the dimension of the exogenous variables  $\mathbf{Z}_t$ ,  $p_n$ , may diverge at an exponential rate of  $n$ , which implies that the potential explanatory variables  $\mathbf{X}_t$  have the dimension of  $p_n + d_n$  diverging at an exponential rate, i.e.,  $p_n + d_n = O(\exp\{n^{\delta_0}\})$  for some positive constant  $\delta_0$ . To ensure that our semiparametric model averaging scheme is feasible both theoretically and numerically, we need to reduce the dimension of the potential covariates  $\mathbf{X}_t$  and select those variables that make a significant contribution to predicting the response. In this paper we propose two schemes to achieve the purpose of dimension reduction. The first scheme is called as the “*KSIS+PMAMAR*” method. It reduces the dimension of the potential covariates by first using the approach of *Kernel Sure Independence Screening* (KSIS), motivated by Fan and Lv (2008), to screen out the unimportant marginal regression (or auto-regression) functions, and then apply the so-called *Penalized Model Averaging MArginal Regression* (PMAMAR) to further select the most relevant regression functions. The second scheme is called the “*PCA+PMAMAR*” method. In this scheme, we assume that the ultra-high dimensional exogenous regressors  $\mathbf{Z}_t$  satisfy an approximate factor model which has been popular in many fields including economics and finance (c.f., Chamberlain and Rothschild, 1983; Fama and French, 1992; Stock and Watson, 2002; Bai and Ng, 2002, 2006), and estimate the factor regressors using the *Principal Component Analysis* (PCA). Then, similarly to the second step in the first scheme, the PMAMAR method is applied to further select the significant estimated factor regressors and auto-regressors.

Under some regularity conditions, we develop the asymptotic properties of the proposed methods. For the KSIS procedure, we establish the sure screening property, indicating that the covariates whose marginal regression functions make a truly significant contribution to estimating the multivariate regression function  $m(\mathbf{x})$  would be selected with probability approaching to one to form a set of the regressors that would undergo a further selection in the PMAMAR procedure. The optimal weight estimation obtained in the PMAMAR procedure is proved to have the well-known sparsity and oracle property that the estimated values of the true zero weights are forced to be zero. For the PCA approach, we show that the estimated latent factors are uniformly consistent at a convergence rate that depends on both  $n$  and  $p_n$ , and the kernel estimation of the marginal regression with the estimated factor regressors is asymptotically equivalent to the same procedure with the rotated true factor regressors. Furthermore, extensions of the proposed semiparametric approaches such as an iterative KSIS+PMAMAR procedure will be discussed. In the simulation studies, we find that our methods outperform some existing methods in terms of forecasting accuracy. We finally apply our methods to forecasting quarterly inflation in the UK.

The rest of the paper is organized as follows. The two semiparametric model averaging schemes

are proposed in Section 2. The asymptotic theory for them is then developed in Section 3. Section 4 discusses some extensions when the methods are implemented in practice. Numerical studies are reported in Section 5 including two simulated examples and one empirical data example. Section 6 concludes. Proofs of the asymptotic results are given in a supplemental document.

## 2 Semiparametric model averaging

In this section, we propose two semiparametric model averaging approaches, which are named as the KSIS+PMAMAR and the PCA+PMAMAR in Sections 2.1 and 2.2, respectively.

### 2.1 *KSIS+PMAMAR method*

As mentioned in Section 1, the KSIS+PMAMAR method is a two-step procedure. We first generalize the Sure Independence Screening (SIS) method introduced by Fan and Lv (2008) to the ultra-high dimensional dynamic time series and general nonparametric setting to screen out covariates whose nonparametric marginal regression functions have low correlations with the response. Then, for the covariates that have survived the screening, we propose a PMAMAR method with first-stage kernel smoothing to further select the exogenous regressors and the lags of the response variable which make significant contribution to estimating the multivariate regression function, and to determine an optimal linear combination of the significant marginal regression and auto-regression functions.

**Step one: KSIS.** For notational simplicity, we let

$$X_{tj} = \begin{cases} Z_{tj}, & j = 1, \dots, p_n, \\ Y_{t-(j-p_n)}, & j = p_n + 1, \dots, p_n + d_n. \end{cases}$$

To measure the contribution made by the univariate covariate  $X_{tj}$  to estimating the multivariate regression function  $m(\mathbf{x}) = \mathbb{E}(Y_t | \mathbf{X}_t = \mathbf{x})$ , we consider the marginal regression function defined by

$$m_j(x_j) = \mathbb{E}(Y_t | X_{tj} = x_j), \quad j = 1, \dots, p_n + d_n,$$

which is the projection of  $Y_t$  onto the univariate component space spanned by  $X_{tj}$ . This function can also be seen as the solution to the following nonparametric optimization problem (c.f., Fan *et al*, 2011):

$$\min_{g_j \in \mathcal{L}_2(\mathbf{P})} \mathbb{E}[Y_t - g_j(X_{tj})]^2,$$

where  $\mathcal{L}_2(\mathbf{P})$  is the class of square integrable functions under the probability measure  $\mathbf{P}$ . We estimate the functions  $m_j(\cdot)$  by the commonly-used kernel smoothing method, although other nonparametric estimation methods such as the local polynomial smoothing and smoothing spline method are also

applicable. The kernel smoother of  $m_j(x_j)$  is

$$\hat{m}_j(x_j) = \frac{\sum_{t=1}^n Y_t K_{tj}(x_j)}{\sum_{t=1}^n K_{tj}(x_j)}, \quad K_{tj}(x_j) = K\left(\frac{X_{tj} - x_j}{h_1}\right), \quad j = 1, \dots, p_n + d_n, \quad (2.1)$$

where  $K(\cdot)$  is a kernel function and  $h_1$  is a bandwidth. To make the above kernel estimation method feasible, we assume that the initial observations,  $Y_0, Y_{-1}, \dots, Y_{-d_n+1}$ , of the response are available.

When the observations are independent and the response variable has zero mean, the paper of [Fan \*et al\* \(2011\)](#) ranks the importance of the covariates by calculating the  $\mathcal{L}_2$ -norm of  $\hat{m}_j(\cdot)$ , and chooses those covariates whose corresponding norms are larger than a pre-determined threshold that usually tends to zero. However, in our time series setting, for  $j$  such that  $j - p_n > 0$ ,  $X_{tj}$  becomes the lag of the response variable  $Y_{t-(j-p_n)}$  and  $m_j(\cdot) = \mathbf{E}(Y_t | Y_{t-(j-p_n)} = \cdot)$ . For time series that are stationary and weakly dependent, it is often reasonable to assume that  $\mathbf{E}(Y_t | Y_{t-(j-p_n)}) \xrightarrow{P} \mathbf{E}(Y_t)$  when  $j - p_n \rightarrow \infty$ . On the other hand, under some regularity conditions, using the uniform consistency result for the kernel smoothing method (c.f., [Li \*et al\*, 2012](#)), we have  $\hat{m}_j(x_j) \xrightarrow{P} m_j(x_j)$  uniformly for  $x_j$  in a compact set. Combining the above arguments, we may show that as  $j - p_n \rightarrow \infty$

$$\hat{m}_j(x_j) \xrightarrow{P} m_j(x_j) \rightarrow \mathbf{E}(Y_t)$$

uniformly for  $x_j$  in a compact set. When  $\mathbf{E}(Y_t)$  is non-zero, the norm of  $\hat{m}_j(\cdot)$  would tend to a non-zero quantity when  $j - p_n \rightarrow \infty$ . As a consequence, if covariates are chosen according to the  $\mathcal{L}_2$ -norm of their corresponding marginal regression functions, quite a few unimportant lags might be chosen. To address this issue, we consider ranking the importance of the covariates by calculating the correlation between the response variable and marginal regression

$$\text{cor}(j) = \frac{\text{cov}(j)}{\sqrt{\mathbf{v}(Y) \cdot \mathbf{v}(j)}} = \left[ \frac{\mathbf{v}(j)}{\mathbf{v}(Y)} \right]^{1/2}, \quad (2.2)$$

where  $\mathbf{v}(Y) = \text{var}(Y_t)$ ,  $\mathbf{v}(j) = \text{var}(m_j(X_{tj}))$  and  $\text{cov}(j) = \text{cov}(Y_t, m_j(X_{tj})) = \text{var}(m_j(X_{tj})) = \mathbf{v}(j)$ . Equation (2.2) indicates that the value of  $\text{cor}(j)$  is non-negative for all  $j$  and the ranking of  $\text{cor}(j)$  is equivalent to the ranking of  $\mathbf{v}(j)$  as  $\mathbf{v}(Y)$  is positive and invariant across  $j$ . The sample version of  $\text{cor}(j)$  can be constructed as

$$\hat{\text{cor}}(j) = \frac{\hat{\text{cov}}(j)}{\sqrt{\hat{\mathbf{v}}(Y) \cdot \hat{\mathbf{v}}(j)}} = \left[ \frac{\hat{\mathbf{v}}(j)}{\hat{\mathbf{v}}(Y)} \right]^{1/2}, \quad (2.3)$$

where:

$$\hat{v}(Y) = \frac{1}{n} \sum_{t=1}^n Y_t^2 - \left( \frac{1}{n} \sum_{t=1}^n Y_t \right)^2,$$

$$\hat{\text{cov}}(j) = \hat{v}(j) = \frac{1}{n} \sum_{t=1}^n \hat{m}_j^2(X_{tj}) - \left[ \frac{1}{n} \sum_{t=1}^n \hat{m}_j(X_{tj}) \right]^2, \quad j = 1, 2, \dots, p_n + d_n.$$

The screened sub-model can be determined by

$$\hat{S} = \{j \in \{1, 2, \dots, p_n + d_n\} : \hat{v}(j) \geq \rho_n\}, \quad (2.4)$$

where  $\rho_n$  is a pre-determined positive number. By (2.3), the criterion in (2.4) is equivalent to

$$\hat{S} = \{j \in \{1, 2, \dots, p_n + d_n\} : \hat{\text{cov}}(j) \geq \rho_n^\diamond\},$$

where  $\rho_n^\diamond = \rho_n^{1/2} / \sqrt{\hat{v}(Y)}$ . We let  $\mathbf{X}_t^* = (X_{t1}^*, X_{t2}^*, \dots, X_{tq_n}^*)^\top$  be the covariates chosen according to the criterion (2.4).

The above model selection procedure can be seen as the nonparametric kernel extension of the SIS method introduced by Fan and Lv (2008) in the context of linear regression models. Recent extensions to nonparametric additive models and varying coefficient models can be found in Fan *et al* (2011), Fan *et al* (2014) and Liu *et al* (2014). However, the existing literature usually considers the case where the observations are either independent or collected from correlated and sparse longitudinal data (c.f., Cheng *et al*, 2014), which rules out the nonlinear dynamic time series setting (over a long time span). In this paper, we relax such a restriction and show that the KSIS approach works well in the ultra-high dimensional time series and semiparametric setting. Also, differently from Fan *et al* (2011) using the B-splines method, our paper applies the kernel smoothing method to estimate the marginal regression functions, with different mathematical tool required to derive our asymptotic theory.

**Step two: PMAMAR.** In the second step, we propose using a semiparametric method of model averaging lower dimensional regression functions to estimate

$$m^*(\mathbf{x}) = \mathbb{E}(Y_t | \mathbf{X}_t^* = \mathbf{x}), \quad (2.5)$$

where  $\mathbf{x} = (x_1, x_2, \dots, x_{q_n})^\top$ . Specifically, we approximate the conditional regression function  $m^*(\mathbf{x})$  by an affine combination of one-dimensional conditional component regressions

$$m_j^*(x_j) = \mathbb{E}(Y_t | X_{tj}^* = x_j), \quad j = 1, \dots, q_n.$$

Each marginal regression  $m_j^*(\cdot)$  can be treated as a “nonlinear candidate model” and the number of such nonlinear candidate models is  $q_n$ . A weighted average of  $m_j^*(x_j)$  is then used to approximate

$m^*(\mathbf{x})$ , i.e.,

$$m^*(\mathbf{x}) \approx w_0 + \sum_{j=1}^{q_n} w_j m_j^*(x_j), \quad (2.6)$$

where  $w_j$ ,  $j = 0, 1, \dots, q_n$ , are to be determined later and can be seen as the weights for different candidate models. The linear combination in (2.6) is called as Model Averaging MARGinal Regressions or MAMAR (c.f., Li *et al*, 2015) and is applied by Chen *et al* (2016) in the dynamic portfolio choice with many conditioning variables. As the conditional component regressions  $m_j^*(X_{tj}^*) = \mathbf{E}(Y_t | X_{tj}^*)$ ,  $j = 1, \dots, q_n$ , are unknown but univariate, in practice, they can be well estimated by various nonparametric approaches that would not suffer from the curse of dimensionality problem. Hence, the first stage in the semiparametric PMAMAR procedure is to estimate the marginal regression functions  $m_j^*(\cdot)$  by the kernel smoothing method

$$\hat{m}_j^*(x_j) = \frac{\sum_{t=1}^n Y_t \bar{K}_{tj}(x_j)}{\sum_{t=1}^n \bar{K}_{tj}(x_j)}, \quad \bar{K}_{tj}(x_j) = K\left(\frac{X_{tj}^* - x_j}{h_2}\right), \quad j = 1, \dots, q_n, \quad (2.7)$$

where  $h_2$  is a bandwidth. Let

$$\hat{\mathcal{M}}(j) = [\hat{m}_j^*(X_{1j}^*), \dots, \hat{m}_j^*(X_{nj}^*)]^\top$$

be the estimated values of

$$\mathcal{M}(j) = [m_j^*(X_{1j}^*), \dots, m_j^*(X_{nj}^*)]^\top$$

for  $j = 1, \dots, q_n$ . By using (2.7), we have

$$\hat{\mathcal{M}}(j) = \mathcal{S}_n(j) \mathcal{Y}_n, \quad j = 1, \dots, q_n,$$

where  $\mathcal{S}_n(j)$  is the  $n \times n$  smoothing matrix whose  $(k, l)$ -component is  $\bar{K}_{lj}(X_{kj}^*) / [\sum_{t=1}^n \bar{K}_{tj}(X_{kj}^*)]$ , and  $\mathcal{Y}_n = (Y_1, \dots, Y_n)^\top$ .

The second stage of PMAMAR is to replace  $m_j^*(X_{tj}^*)$ ,  $j = 1, \dots, q_n$ , by their corresponding nonparametric estimates  $\hat{m}_j^*(X_{tj}^*)$ , and use the penalized approach to select the significant marginal regression functions in the following ‘‘approximate linear model’’:

$$Y_t \approx w_0 + \sum_{j=1}^{q_n} w_j \hat{m}_j^*(X_{tj}^*). \quad (2.8)$$

Without loss of generality, we further assume that  $\mathbf{E}(Y_t) = 0$ , otherwise, we can simply replace  $Y_t$  by  $Y_t - \bar{Y} = Y_t - \frac{1}{n} \sum_{s=1}^n Y_s$ . It is easy to show that the intercept term  $w_0$  in (2.6) is zero under this assumption. In the sequel, we let  $\mathbf{w}_o := \mathbf{w}_{o_n} = (w_{o1}, \dots, w_{oq_n})$  be the optimal values of the weights in the model averaging defined as in Li *et al* (2015). Based on the approximate linear modelling



framework (2.8), for given  $\mathbf{w}_n = (w_1, \dots, w_{q_n})^\top$ , we define the objective function by

$$\mathcal{Q}_n(\mathbf{w}_n) = [\mathcal{Y}_n - \hat{\mathcal{M}}(\mathbf{w}_n)]^\top [\mathcal{Y}_n - \hat{\mathcal{M}}(\mathbf{w}_n)] + n \sum_{j=1}^{q_n} p_\lambda(|w_j|), \quad (2.9)$$

where

$$\hat{\mathcal{M}}(\mathbf{w}_n) = [w_1 \mathcal{S}_n(1) + \dots + w_{q_n} \mathcal{S}_n(q_n)] \mathcal{Y}_n = \mathcal{S}_n(\mathcal{Y}) \mathbf{w}_n,$$

$\mathcal{S}_n(\mathcal{Y}) = [\mathcal{S}_n(1) \mathcal{Y}_n, \dots, \mathcal{S}_n(q_n) \mathcal{Y}_n]$ , and  $p_\lambda(\cdot)$  is a penalty function with a tuning parameter  $\lambda$ . The vector  $\hat{\mathcal{M}}(\mathbf{w}_n)$  in (2.6) can be seen as the kernel estimate of

$$\mathcal{M}(\mathbf{w}_n) = \left[ \sum_{j=1}^{q_n} w_j m_j^*(X_{1j}^*), \dots, \sum_{j=1}^{q_n} w_j m_j^*(X_{nj}^*) \right]^\top$$

for given  $\mathbf{w}_n$ . Our semiparametric estimator of the optimal weights  $\mathbf{w}_o$  can be obtained through minimizing the objective function  $\mathcal{Q}_n(\mathbf{w}_n)$ :

$$\hat{\mathbf{w}}_n = \arg \min_{\mathbf{w}_n} \mathcal{Q}_n(\mathbf{w}_n). \quad (2.10)$$

There has been extensive discussion on the choice of the penalty function for parametric linear and nonlinear models. Many popular variable selection criteria, such as AIC and BIC, correspond to the penalized estimation method with  $p_\lambda(|z|) = 0.5\lambda^2 I(|z| \neq 0)$  with different values of  $\lambda$ . However, as mentioned by Fan and Li (2001), such traditional penalized approaches are expensive in computational cost when  $q_n$  is large. To avoid the computational burden and the lack of stability, some other penalty functions have been introduced in recent years. For example, the LASSO penalty  $p_\lambda(|z|) = \lambda|z|$  has been extensively studied by many authors (c.f., Tibshirani, 1996, 1997); Frank and Friedman (1993) consider the  $L_q$ -penalty  $p_\lambda(|z|) = \lambda|z|^q$  for  $0 < q < 1$ ; Fan and Li (2001) suggest using the SCAD penalty function whose derivative is defined by

$$p'_\lambda(z) = \lambda \left[ I(z \leq \lambda) + \frac{a_0 \lambda - z}{(a_0 - 1)\lambda} I(z > \lambda) \right]$$

with  $p_\lambda(0) = 0$ , where  $a_0 > 2$ ,  $\lambda > 0$  and  $I(\cdot)$  is the indicator function.

## 2.2 PCA+PMAMAR method

Because of dependence within the exogenous variables in  $\mathbf{Z}_t$ , sparsity may be an issue when its dimension is high. It is well known that we may also achieve dimension reduction through the use of factor models when analyzing high-dimensional time series data. In this subsection, we assume that

the high-dimensional exogenous variables  $\mathbf{Z}_t$  follow the approximate factor model:

$$Z_{tk} = (\mathbf{b}_k^0)^\top \mathbf{f}_t^0 + u_{tk}, \quad k = 1, \dots, p_n, \quad (2.11)$$

where  $\mathbf{b}_k^0$  is an  $r$ -dimensional vector of factor loadings,  $\mathbf{f}_t^0$  is an  $r$ -dimensional vector of common factors, and  $u_{tk}$  is called an idiosyncratic error. The number of the common factors,  $r$ , is assumed to be fixed throughout the paper, but it is usually unknown in practice and its determination method will be discussed in Section 4 below.

From the approximate factor model (2.11), we can find that the main information in the exogenous regressors may be summarized in the common factors  $\mathbf{f}_t^0$  that have a much lower dimension. The aim of dimension reduction can thus be achieved, and it may be reasonable to replace  $\mathbf{Z}_t$  with an ultra-high dimension by the unobservable  $\mathbf{f}_t$  with a fixed dimension in estimating the conditional multivariate regression function and predicting the future value of the response variable  $Y_t$ . In the framework of linear regression or autoregression, such an idea has been frequently used in the literature since [Stock and Watson \(2002\)](#) and [Bernanke \*et al\* \(2005\)](#). However, so far as we know, there is virtually no work on combining the factor model (2.11) with the nonparametric nonlinear regression. The only exception is the paper by [Härdle and Tsybakov \(1995\)](#), which considers the additive regression model on principal components when the observations are independent and the dimension of the potential regressors is fixed. The latter restriction is relaxed in this paper.

Instead of directly studying the multivariate regression function  $m(\mathbf{x})$  defined in (1.1), we next consider the multivariate regression function defined by

$$m_f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{E} (Y_t | \mathbf{f}_t^0 = \mathbf{x}_1, \mathbf{Y}_{t-1} = \mathbf{x}_2), \quad (2.12)$$

where  $\mathbf{Y}_{t-1}$  is defined as in Section 1,  $\mathbf{x}_1$  is  $r$ -dimensional and  $\mathbf{x}_2$  is  $d_n$ -dimensional. In order to develop a feasible estimation approach for the factor augmented nonlinear regression function in (2.12), we need to estimate the unobservable factor regressors  $\mathbf{f}_t^0$  in the first step. This will be done through the PCA approach and we denote

$$\hat{\mathbf{X}}_{t,f}^* = \left( \hat{\mathbf{f}}_t^\top, \mathbf{Y}_{t-1}^\top \right)^\top = \left( \hat{f}_{t1}, \dots, \hat{f}_{tr}, \dots, \mathbf{Y}_{t-1}^\top \right)^\top$$

as a combination of the estimated factor regressors and lags of response variables, where  $\hat{\mathbf{f}}_t$  is the estimated factor via PCA and  $\hat{f}_{tk}$  is the  $k$ -th element of  $\hat{\mathbf{f}}_t$ ,  $k = 1, \dots, r$ . In the second step, we use the PMAMAR method to conduct a further selection among the  $(r + d_n)$ -dimensional regressors and determine an optimal combination of the significant marginal regressions. This PCA+PMAMAR method substantially generalizes the framework of factor-augmented linear regression or autoregression (c.f., [Stock and Watson, 2002](#); [Bernanke \*et al\*, 2005](#); [Bai and Ng, 2006](#); [Pesaran \*et al\*, 2011](#); [Cheng and Hansen, 2015](#)) to the general semiparametric framework.

**Step one: PCA on the exogenous regressors.** Letting

$$\mathbf{B}_n^0 = (\mathbf{b}_1^0, \dots, \mathbf{b}_{p_n}^0)^\top \quad \text{and} \quad \mathbf{U}_t = (u_{t1}, \dots, u_{tp_n})^\top,$$

we may rewrite the approximate factor model (2.11) as

$$\mathbf{Z}_t = \mathbf{B}_n^0 \mathbf{f}_t^0 + \mathbf{U}_t. \quad (2.13)$$

We next apply the PCA approach to obtain the estimation of the common factors  $\mathbf{f}_t^0$ . Denote  $\mathcal{Z}_n = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)^\top$ , the  $n \times p_n$  matrix of the observations of the exogenous variables. We then construct  $\hat{\mathcal{F}}_n = (\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_n)^\top$  as the  $n \times r$  matrix consisting of the  $r$  eigenvectors (multiplied by  $\sqrt{n}$ ) associated with the  $r$  largest eigenvalues of the  $n \times n$  matrix  $\mathcal{Z}_n \mathcal{Z}_n^\top / (np_n)$ . Furthermore, the estimate of the factor loading matrix (with rotation) is defined as

$$\hat{\mathbf{B}}_n = (\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_{p_n})^\top = \mathcal{Z}_n^\top \hat{\mathcal{F}}_n / n,$$

by noting that  $\hat{\mathcal{F}}_n^\top \hat{\mathcal{F}}_n / n = I_r$ .

As shown in the literature (see also Theorem 3 in Section 3.2 below),  $\hat{\mathbf{f}}_t$  is a consistent estimator of the rotated common factor  $\mathbf{H} \mathbf{f}_t^0$ , where

$$\mathbf{H} = \hat{\mathbf{V}}^{-1} \left( \hat{\mathcal{F}}_n^\top \mathcal{F}_n^0 / n \right) [(\mathbf{B}_n^0)^\top \mathbf{B}_n^0 / p_n], \quad \mathcal{F}_n^0 = (\mathbf{f}_1^0, \dots, \mathbf{f}_n^0)^\top,$$

and  $\hat{\mathbf{V}}$  is the  $r \times r$  diagonal matrix of the first  $r$  largest eigenvalues of  $\mathcal{Z}_n \mathcal{Z}_n^\top / (np_n)$  arranged in descending order. Consequently, we may consider the following multivariate regression function with rotated latent factors:

$$m_f^*(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{E} (Y_t | \mathbf{H} \mathbf{f}_t^0 = \mathbf{x}_1, \mathbf{Y}_{t-1} = \mathbf{x}_2). \quad (2.14)$$

In the subsequent PMAMAR step, we can use  $\hat{\mathbf{f}}_t$  to replace  $\mathbf{H} \mathbf{f}_t^0$  in the semiparametric procedure. The factor modelling and PCA estimation ensure that most of the useful information contained in the exogenous variables  $\mathbf{Z}_t$  can be extracted before the second step of PMAMAR, which may lead to possible good performance in forecasting  $Y_t$  through the use of the estimated common factors. In contrast, as discussed in some existing literature such as Fan and Lv (2008), when irrelevant exogenous variables are highly correlated with some relevant ones, they might be selected into a model by the SIS or KSIS procedure with higher priority than some other relevant exogenous variables, which results in high false positive rates and low true positive rates and leads to loss of useful information in the potential covariates, see, for example, the discussion in Section 4.1.

**Step two: PMAMAR using estimated factor regressors.** Recall that

$$\hat{\mathbf{X}}_{t,f}^* = \left( \hat{\mathbf{f}}_t^\top, \mathbf{Y}_{t-1}^\top \right)^\top = \left( \hat{f}_{t1}, \dots, \hat{f}_{tr}, \mathbf{Y}_{t-1}^\top \right)^\top,$$

where  $\hat{f}_{tk}$  is the  $k$ -th element of  $\hat{\mathbf{f}}_t$ ,  $k = 1, \dots, r$ . We may apply the two-stage semiparametric PMAMAR procedure, which is exactly the same as that in Section 2.1 to the process  $(Y_t, \hat{\mathbf{X}}_{t,f}^*)$ ,  $t = 1, \dots, n$ , and then obtain the estimation of the optimal weights  $\hat{\mathbf{w}}_{n,f}$ . To save space, we next only sketch the kernel estimation of the marginal regression function with the estimated factor regressors obtained via PCA.

For  $k = 1, \dots, r$ , define

$$m_{k,f}^*(z_k) = \mathbb{E} \left[ Y_t | \tilde{f}_{tk}^0 = z_k \right], \quad \tilde{f}_{tk}^0 = e_r^\top(k) \mathbf{H} \mathbf{f}_t^0,$$

where  $e_r(k)$  is an  $r$ -dimensional column vector with the  $k$ -th element being one and zeros elsewhere,  $k = 1, \dots, r$ . As in Section 2.1, we estimate  $m_{k,f}^*(z_k)$  by the kernel smoothing method:

$$\hat{m}_{k,f}^*(z_k) = \frac{\sum_{t=1}^n Y_t \tilde{K}_{tk}(z_k)}{\sum_{t=1}^n \tilde{K}_{tk}(z_k)}, \quad \tilde{K}_{tk}(z_k) = K\left(\frac{\hat{f}_{tk} - z_k}{h_3}\right), \quad j = 1, \dots, r, \quad (2.15)$$

where  $h_3$  is a bandwidth. In Section 3.2 below, we will show that  $\hat{m}_{k,f}^*(z_k)$  is asymptotically equivalent to  $\tilde{m}_{k,f}^*(z_k)$ , which is defined as in (2.15) but with  $\hat{f}_{tk}$  replaced by  $\tilde{f}_{tk}^0$ . The latter kernel estimation is infeasible in practice as the factor regressor involved is unobservable. As we may show that the asymptotic order of  $\hat{m}_{k,f}^*(z_k) - \tilde{m}_{k,f}^*(z_k)$  is  $o_P(n^{-1/2})$  under some mild conditions (c.f., Theorem 3), the influence of replacing  $\tilde{f}_{tk}^0$  by the estimated factor regressors  $\hat{f}_{tk}$  in the PMAMAR procedure is asymptotically negligible.

### 3 The main theoretical results

In this section, we establish the asymptotic properties for the methodologies developed in Section 2 above. The asymptotic theory for the KSIS+PMAMAR method is given in Section 3.1 and that for the PCA+PMAMAR method is given in Section 3.2.

#### 3.1 Asymptotic theory for KSIS+PMAMAR

In this subsection, we first derive the sure screening property for the developed KSIS method, which implies that the covariates whose marginal regression functions make significant contribution to estimating the multivariate regression function  $m(\mathbf{x})$  would be chosen in the screening with probability approaching one. The following regularity conditions are needed in the proof of this property.

**A1.** *The process  $\{(Y_t, \mathbf{X}_t)\}$  is stationary and  $\alpha$ -mixing with the mixing coefficient decaying at a geometric rate:  $\alpha(k) \sim c_\alpha \theta_0^k$ , where  $0 < c_\alpha < \infty$  and  $0 < \theta_0 < 1$ .*

**A2.** *Let  $f_j(\cdot)$  be the marginal density function of  $X_{tj}$ , the  $j$ -th element of  $\mathbf{X}_t$ . Assume that  $f_j(\cdot)$*

has continuous derivatives up to the second order and

$$0 < \underline{c} \leq \inf_j \inf_{x_j \in \mathcal{C}_j} f_j(x_j) \leq \sup_j \sup_{x_j \in \mathcal{C}_j} f_j(x_j) \leq \bar{c} < \infty,$$

where  $\mathcal{C}_j$  is the compact support of  $X_{tj}$ . For each  $j$ , the conditional density functions of  $Y_t$  for given  $X_{tj}$  exists and satisfies the Lipschitz continuous condition. Furthermore, the length of  $\mathcal{C}_j$  is uniformly bounded by a positive constant.

- A3.** The kernel function  $K(\cdot)$  is a Lipschitz continuous, symmetric and bounded probability density function with a compact support. Let the bandwidth satisfy  $h_1 \sim n^{-\theta_1}$  with  $1/6 < \theta_1 < 1$ .
- A4.** The marginal regression function  $m_j(\cdot)$  has continuous derivatives up to the second order and there exists a positive constant  $c_m$  such that  $\sup_j \sup_{x_j \in \mathcal{C}_j} [ |m_j(x_j)| + |m'_j(x_j)| + |m''_j(x_j)| ] \leq c_m$ .
- A5.** The response variable  $Y_t$  satisfies  $E[\exp\{\varsigma|Y_t|\}] < \infty$ , where  $\varsigma$  is a positive constant.

**Remark 1.** The condition **A1** imposes the stationary  $\alpha$ -mixing dependence structure on the observations, which is not uncommon in the time series literature (c.f., [Bosq, 1998](#)). It might be possible to consider a more general dependence structure such as the near epoch dependence studied in [Lu and Linton \(2007\)](#) and [Li et al \(2012\)](#), however, the technical proofs would be more involved. Hence, we impose the mixing dependence structure and focus on the ideas proposed. The restriction of geometric decaying rate on the mixing coefficient is due to the ultra-high dimensional setting and it may be relaxed if the dimension of the covariates diverges at a polynomial rate. The conditions **A2** and **A4** give some smoothness restrictions on the marginal density functions and marginal regression functions. To simplify the discussion, we assume that all of the marginal density functions have compact support. Such an assumption might be too restrictive for time series data, but it could be relaxed by slightly modifying our methodology. For example, if the marginal density function of  $X_{tj}$  is the standard normal density which does not have a compact support, we can truncate the tail of  $X_{tj}$  in the KSIS procedure by replacing  $X_{tj}$  with  $X_{tj}I(|X_{tj}| \leq \zeta_n)$  and  $\zeta_n$  divergent to infinity at a slow rate. The condition **A3** is a commonly-used condition on the kernel function as well as the bandwidth. The strong moment condition on  $Y_t$  in **A5** is also quite common in the SIS literature such as [Fan et al \(2011\)](#) and [Liu et al \(2014\)](#).

Define the index set of “true” candidate models as

$$\mathcal{S} = \{j = 1, 2, \dots, p_n + d_n : v(j) \neq 0\}.$$

The following theorem gives the sure screening property for the KSIS procedure.

**Theorem 1.** *Suppose that the conditions A1–A5 are satisfied.*

(i) For any small  $\delta_1 > 0$ , there exists a positive constant  $\delta_2$  such that

$$\mathbb{P} \left( \max_{1 \leq j \leq p_n + d_n} \left| \hat{\mathbf{v}}(j) - \mathbf{v}(j) \right| > \delta_1 n^{-2(1-\theta_1)/5} \right) = O \left( M(n) \exp \left\{ -\delta_2 n^{(1-\theta_1)/5} \right\} \right), \quad (3.1)$$

where  $M(n) = (p_n + d_n)n^{(17+18\theta_1)/10}$  and  $\theta_1$  is defined in the condition A3.

(ii) If we choose the pre-determined tuning parameter  $\rho_n = \delta_1 n^{-2(1-\theta_1)/5}$  and assume

$$\min_{j \in \mathcal{S}} \mathbf{v}(j) \geq 2\delta_1 n^{-2(1-\theta_1)/5}, \quad (3.2)$$

then we have

$$\mathbb{P}(\mathcal{S} \subset \hat{\mathcal{S}}) \geq 1 - O \left( M_{\mathcal{S}}(n) \exp \left\{ -\delta_2 n^{(1-\theta_1)/5} \right\} \right), \quad (3.3)$$

where  $M_{\mathcal{S}}(n) = |\mathcal{S}|n^{(17+18\theta_1)/10}$  with  $|\mathcal{S}|$  being the cardinality of  $\mathcal{S}$ .

**Remark 2.** The above theorem shows that the covariates whose marginal regressions have not too small positive correlations with the response variable would be included in the screened model with probability approaching one at a possible exponential rate of  $n$ . The condition (3.2) guarantees that the correlations between the response and the marginal regression functions for covariates whose indices belong to  $\mathcal{S}$  are bounded away from zero, but the lower bound may converge to zero. As  $p_n + d_n = O(\exp\{n^{\delta_0}\})$ , in order to ensure the validity of Theorem 1(i), we need to impose the restriction  $\delta_0 < (1 - \theta_1)/5$ , which reduces to  $\delta_0 < 4/25$  if the order of the optimal bandwidth in kernel smoothing (i.e.,  $\theta_1 = 1/5$ ) is used. Our theorem generalizes the results in Fan *et al* (2011) and Liu *et al* (2014) to dynamic time series case and those in Ando and Li (2014) to the flexible nonparametric setting.

We next study the asymptotic properties for the PMAMAR method including the well-known sparsity and oracle property. Recall that  $q_n = |\hat{\mathcal{S}}|$  and the dimension of the potential covariates is reduced from  $p_n + d_n$  to  $q_n$  after implementing the KSIS procedure. As above, we let  $\mathbf{X}_t^*$  be the KSIS-chosen covariates, which may include both the exogenous regressors and lags of  $Y_t$ . Define

$$a_n = \max_{1 \leq j \leq q_n} \left\{ |p'_\lambda(|w_{oj}|)|, |w_{oj}| \neq 0 \right\}$$

and

$$b_n = \max_{1 \leq j \leq q_n} \left\{ |p''_\lambda(|w_{oj}|)|, |w_{oj}| \neq 0 \right\}.$$

We need to introduce some additional conditions to derive the asymptotic theory.

#### A6. The matrix

$$\mathbf{\Lambda}_n := \begin{pmatrix} \mathbb{E}[m_1^*(X_{t1}^*)m_1^*(X_{t1}^*)] & \dots & \mathbb{E}[m_1^*(X_{t1}^*)m_{q_n}^*(X_{tq_n}^*)] \\ \vdots & \vdots & \vdots \\ \mathbb{E}[m_{q_n}^*(X_{tq_n}^*)m_1^*(X_{t1}^*)] & \dots & \mathbb{E}[m_{q_n}^*(X_{tq_n}^*)m_{q_n}^*(X_{tq_n}^*)] \end{pmatrix}$$

is positive definite with the eigenvalues bounded away from zero and infinity. In particular, the smallest eigenvalue of  $\Lambda_n$  is larger than  $\chi$ , a small positive constant.

**A7.** The bandwidth  $h_2$  satisfies

$$s_n^2 n h_2^4 \rightarrow 0, \quad n^{\frac{1}{2}-\xi} h_2 \rightarrow \infty, \quad q_n^2 (\tau_n + h_2^2) = o(1) \quad (3.4)$$

as  $n \rightarrow \infty$ , where  $s_n$  is the number of non-zero elements in the optimal weight vector,  $\xi$  is positive but arbitrarily small, and  $\tau_n = \left(\frac{\log n}{n h_2}\right)^{1/2}$ .

**A8.** Let  $a_n = O(n^{-1/2})$ ,  $b_n = o(1)$ ,  $p_\lambda(0) = 0$ , and there exist two positive constants  $C_1$  and  $C_2$  such that  $|p''_\lambda(w_1) - p''_\lambda(w_2)| \leq C_2 |w_1 - w_2|$  when  $w_1, w_2 > C_1 \lambda$ .

**Remark 3.** The condition **A6** gives some regularity conditions on the eigenvalues of the  $q_n \times q_n$  positive definite matrix  $\Lambda_n$ , which are similar to those in the existing literature dealing with independent observations (c.f., [Fan and Peng, 2004](#)). We may relax these conditions by allowing that some eigenvalues tend to zero at certain rates and slightly modifying the conditions **A7** and **A8**, and as a consequence, the convergence rate in Theorem 2(i) below would be slightly different (c.f., [Chen et al, 2015](#)). The restrictions in the condition **A7** imply that undersmoothing is needed in our semiparametric procedure and  $q_n$  can only be divergent at a polynomial rate of  $n$ . The condition **A8** is a commonly-used condition on the penalty function  $p_\lambda(\cdot)$ , similar to that in [Fan and Peng \(2004\)](#).

Without loss of generality, we define the vector of the optimal weights

$$\mathbf{w}_o = (w_{o1}, \dots, w_{oq_n})^\top = [\mathbf{w}_o^\top(1), \mathbf{w}_o^\top(2)]^\top,$$

where  $\mathbf{w}_o(1)$  is composed of non-zero weights with dimension  $s_n$  and  $\mathbf{w}_o(2)$  is composed of zero weights with dimension  $(q_n - s_n)$ , and assume that the observations  $X_{tj}^*$  are in the interior of the respective support (which is to avoid the kernel boundary effect in the asymptotic analysis). In order to give the asymptotic normality for  $\hat{\mathbf{w}}_n(1)$ , the estimator of  $\mathbf{w}_o(1)$ , we need to introduce some further notation. Define

$$\eta_t^* = Y_t - \sum_{j=1}^{q_n} w_{oj} m_j^*(X_{tj}^*), \quad \eta_{tj}^* = Y_t - m_j^*(X_{tj}^*)$$

and  $\boldsymbol{\xi}_t = (\xi_{t1}, \dots, \xi_{ts_n})^\top$  with  $\xi_{tj} = \bar{\eta}_{tj}^* - \tilde{\eta}_{tj}^*$ ,  $\bar{\eta}_{tj}^* = m_j^*(X_{tj}^*) \eta_t^*$ ,

$$\tilde{\eta}_{tj}^* = \sum_{k=1}^{q_n} w_{ok} \eta_{tk}^* \beta_{jk}(X_{tk}^*) = \sum_{k=1}^{s_n} w_{ok} \eta_{tk}^* \beta_{jk}(X_{tk}^*), \quad \beta_{jk}(x_k) = \mathbf{E} [m_j^*(X_{tj}^*) | X_{tk}^* = x_k].$$

Obviously, the mean of  $\boldsymbol{\xi}_t$  is zero, and we define  $\Sigma_n = \sum_{t=-\infty}^{\infty} \mathbf{E}(\boldsymbol{\xi}_0 \boldsymbol{\xi}_t^\top)$  and  $\Lambda_{n1}$  as the top-left  $s_n \times s_n$  submatrix of  $\Lambda_n$ . Let

$$\boldsymbol{\omega}_n = [p'_\lambda(|w_{o1}|) \text{sgn}(w_{o1}), \dots, p'_\lambda(|w_{os_n}|) \text{sgn}(w_{os_n})]^\top$$

and

$$\mathbf{\Omega}_n = \text{diag} \{p''_\lambda(|w_{o1}|), \dots, p''_\lambda(|w_{os_n}|)\},$$

where  $\text{sgn}(\cdot)$  is the sign function. In the following theorem, we give the asymptotic theory of  $\hat{\mathbf{w}}_n$  obtained by the PMAMAR method.

**Theorem 2.** *Suppose that the conditions A1–A8 are satisfied.*

(i) *There exists a local minimizer  $\hat{\mathbf{w}}_n$  of the objective function  $\mathcal{Q}_n(\cdot)$  defined in (2.9) such that*

$$\|\hat{\mathbf{w}}_n - \mathbf{w}_o\| = O_P\left(\sqrt{q_n}(n^{-1/2} + a_n)\right) = O_P\left(\sqrt{q_n/n}\right), \quad (3.5)$$

where  $a_n$  is defined in the condition A8 and  $\|\cdot\|$  denotes the Euclidean norm.

(ii) *Let  $\hat{\mathbf{w}}_n(2)$  be the estimator of  $\mathbf{w}_o(2)$  and further assume that*

$$\lambda \rightarrow 0, \quad \frac{\sqrt{n}\lambda}{\sqrt{q_n}} \rightarrow \infty, \quad \liminf_{n \rightarrow \infty} \liminf_{w \rightarrow 0^+} \frac{p'_\lambda(w)}{\lambda} > 0. \quad (3.6)$$

Then, the local minimizer  $\hat{\mathbf{w}}_n$  of the objective function  $\mathcal{Q}_n(\cdot)$  satisfies  $\hat{\mathbf{w}}_n(2) = \mathbf{0}$  with probability approaching one.

(iii) *Letting  $\hat{\mathbf{w}}_n(1)$  be the estimator of  $\mathbf{w}_o(1)$ ,*

$$\sqrt{n}\mathbf{A}_n\mathbf{\Sigma}_n^{-1/2}(\mathbf{\Lambda}_{n1} + \mathbf{\Omega}_n) \left[ \hat{\mathbf{w}}_n(1) - \mathbf{w}_o(1) - (\mathbf{\Lambda}_{n1} + \mathbf{\Omega}_n)^{-1}\boldsymbol{\omega}_n \right] \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{A}_0), \quad (3.7)$$

where  $\mathbf{0}$  is a null vector whose dimension may change from line to line,  $\mathbf{A}_n$  is an  $s \times s_n$  matrix such that  $\mathbb{E} \left[ \|\mathbf{A}_n\mathbf{\Sigma}_n^{-1/2}\boldsymbol{\xi}_t\|^{2+\delta_*} \right] < \infty$  for some  $\delta_* > 0$  and  $\mathbf{A}_n\mathbf{A}_n^\top \rightarrow \mathbf{A}_0$  in which  $\mathbf{A}_0$  is an  $s \times s$  symmetric and non-negative definite matrix and  $s$  is a fixed positive integer.

**Remark 4.** Theorem 2(i) shows that the convergence rate of the estimator  $\hat{\mathbf{w}}_n$  is the same as that in Theorem 1 of Fan and Peng (2004) who consider the case of independent observations. Furthermore, when  $q_n$  is fixed and  $a_n = O(n^{-1/2})$ , we could derive the root- $n$  convergence rate for  $\hat{\mathbf{w}}_n$  as in Theorem 3.1 of Li et al (2015). Theorem 2(ii) shows that the estimator of  $\mathbf{w}_o(2)$  is equal to zero with probability approaching one, which indicates that the PMAMAR procedure possesses the well known sparsity property, and thus can be used as a model selector. Theorem 2(ii) and (iii) above shows that the proposed estimator of the optimal weights enjoy the oracle property, which takes  $\mathbf{w}_o(2) = \mathbf{0}$  as a prerequisite. Furthermore, when  $n$  is large enough and  $\lambda$  tends to zero sufficiently fast for some penalty functions (such as the SCAD penalty), the asymptotic distribution in (3.7) would reduce to

$$\sqrt{n}\mathbf{A}_n\mathbf{\Sigma}_n^{-1/2}\mathbf{\Lambda}_{n1} \left[ \hat{\mathbf{w}}_n(1) - \mathbf{w}_o(1) \right] \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{A}_0), \quad (3.8)$$

which is exactly the same as that in Theorem 3.3 of Li et al (2015).



### 3.2 Asymptotic theory for PCA+PMAMAR

In this subsection, we show that the estimated common factors consistently estimate the true common factors (with rotation), and the asymptotic order of the difference between  $\hat{m}_{k,f}^*(z_k)$  defined in (2.15) and the infeasible kernel estimation  $\tilde{m}_{k,f}^*(z_k)$  is  $o_P(n^{-1/2})$  uniformly. We start with some regularity conditions that are used when proving the asymptotic results.

- B1.** Given the rotation matrix  $\mathbf{H}$ , the process  $\{(Y_t, \mathbf{f}_t, \mathbf{U}_t) : t = 1, \dots, n\}$  is stationary and  $\alpha$ -mixing with the mixing coefficient decaying at a geometric rate.
- B2.** The random common factors satisfy the conditions that  $\mathbf{E}(\mathbf{f}_t^0) = \mathbf{0}$ ,  $\max_t \|\mathbf{f}_t^0\| = O_P(1)$ , the  $r \times r$  matrix  $\mathbf{\Lambda}_F := \mathbf{E}[\mathbf{f}_t^0(\mathbf{f}_t^0)^\top]$  is positive definite and  $\mathbf{E}[\|\mathbf{f}_t^0\|^{4+\tau}] < \infty$  for some  $0 < \tau < \infty$ .
- B3.** The matrix  $(\mathbf{B}_n^0)^\top \mathbf{B}_n^0 / p_n$  is positive definite with the smallest eigenvalue bounded away from zero and  $\max_k \|\mathbf{b}_k^0\|$  is bounded.
- B4.** The idiosyncratic error satisfies  $\mathbf{E}(u_{tk}) = 0$ ,  $\mathbf{E}(u_{tk}\mathbf{f}_t^0) = \mathbf{0}$  and  $\max_k \mathbf{E}[|u_{tk}|^{16}] < \infty$ . Furthermore, there exist two positive constants  $C_3$  and  $C_4$  such that

$$\max_t \mathbf{E} \left[ \left\| \sum_{k=1}^{p_n} u_{tk} \mathbf{b}_k^0 \right\|^4 \right] \leq C_3 p_n^2 \quad (3.9)$$

and

$$\max_{t_1, t_2} \mathbf{E} \left[ \left| \sum_{k=1}^{p_n} \{u_{t_1 k} u_{t_2 k} - \mathbf{E}[u_{t_1 k} u_{t_2 k}]\} \right|^8 \right] \leq C_4 p_n^4, \quad (3.10)$$

and  $\max_k \mathbf{E}[\exp\{\varsigma \|u_{tk}\mathbf{f}_t^0\|\}] < \infty$  where  $\varsigma$  is a positive constant as in the condition A5.

- B5. (i)** The kernel function  $K(\cdot)$  is positive, symmetric and has continuous derivatives up to the second order with a compact support. In addition, the derivative functions of  $K(\cdot)$  are bounded.
- (ii)** There exists  $1/7 < \gamma_0 < 1/6$  such that  $n^{1-\gamma_0} h_3^3 \rightarrow \infty$ . In addition,  $nh_3^4 = O(1)$ ,  $p_n h_3^9 \rightarrow \infty$ , and  $n = o(p_n^2 h_3^{13})$ .
- (iii)** The marginal regression functions (corresponding to the factor regressors)  $m_{k,f}^*(\cdot)$  have continuous and bounded derivatives up to the second order.

**Remark 5.** Some of the above conditions have been commonly used in the literature. For example, the conditions **B2** and **B3** are similar to Assumptions A and B in Bai and Ng (2002), whereas the condition **B4** is similar to the corresponding conditions in Assumption 3.4 in Fan et al (2013). In particular, the exponential bound  $\max_k \mathbf{E}[\exp\{\varsigma \|u_{tk}\mathbf{f}_t^0\|\}] < \infty$  in the condition **B4** is crucial to ensure that  $p_n$  can diverge at an exponential rate of  $n$ . The conditional mixing condition in **B1** seems somehow restrictive, but may be replaced by some weaker (and high-level) conditions. The technical

conditions in **B5**(ii) indicate that the dimension  $p_n$  diverges to infinity at a faster rate than the time series length  $n$  (a commonly-used setting in high-dimensional factor analysis), which are mainly used for the proof of Theorem 3(ii) in Appendix B.

**Theorem 3.** *Suppose that the conditions B1–B4 are satisfied, and*

$$n = o(p_n^2), \quad p_n = O(\exp\{n^{\delta_*}\}), \quad 0 \leq \delta_* < 1/3. \quad (3.11)$$

(i) *For the PCA estimation  $\hat{\mathbf{f}}_t$ , we have*

$$\max_t \left\| \hat{\mathbf{f}}_t - \mathbf{H}\mathbf{f}_t^0 \right\| = O_P(n^{-1/2} + n^{1/4}p_n^{-1/2}), \quad (3.12)$$

where  $\mathbf{H}$  is defined in Section 2.2.

(ii) *In addition, suppose that the conditions A5 and B5 are satisfied and the latent factor  $\mathbf{f}_t^0$  has a compact support. Then we have*

$$\max_{1 \leq k \leq r} \sup_{z_k \in \mathcal{F}_k^*} \left| \hat{m}_{k,f}^*(z_k) - \tilde{m}_{k,f}^*(z_k) \right| = o_P(n^{-1/2}), \quad (3.13)$$

where  $\mathcal{F}_k^*$  is the compact support of  $f_{tk}^0$ .

**Remark 6.** Theorem 3(i) gives the uniform consistency result for the estimation of the common factors, which is very similar to some existing results on PCA estimation of the high-dimensional factor models such as Theorem 3.3 in Fan *et al* (2013). If we further assume that  $n^3 = o(p_n^2)$ , which automatically holds when  $p_n$  is divergent at an exponential rate of  $n$ , the uniform convergence rate in (3.12) would be  $O_P(n^{-1/2})$ . Theorem 3(ii) shows that we may replace  $\hat{m}_{k,f}^*(\cdot)$  by the infeasible kernel estimation  $\tilde{m}_{k,f}^*(\cdot)$  when deriving the asymptotic theory for the PMAMAR method introduced in Section 3.2, and Theorem 2 in Section 3.1 may hold with some notational modifications (c.f.,  $q_n$  in (3.5) needs to be replaced by  $d_n$ ). The restriction of compact support on  $\mathbf{f}_t^0$  can be removed if we slightly modify the methodology as discussed in Remark 1.

## 4 Some extensions

This section discusses some extensions by introducing an iterative KSIS+PMAMAR procedure when the covariates are highly correlated, and an extended PCA+PMAMAR approach with selection of the number of the latent factors in model (2.11).

## 4.1 *An iterative KSIS+PMAMAR procedure*

When the covariates are highly correlated with each other, difficulties in variable selection arise. As documented in [Fan and Lv \(2008\)](#), when the covariate dimension is large, even if the covariates are mutually independent, the data generated from them may exhibit significant spurious correlation. [Fan and Lv \(2008\)](#) notice that when irrelevant covariates are highly correlated with some relevant ones, they might be selected into a model with higher priority than some other relevant covariates, which results in high false positive rates and low true positive rates. Such a problem may become even worse in this paper due to the time series nature of the data, where both the response  $Y_t$  and the covariates  $\mathbf{X}_t$  are autocorrelated over time  $t$ . Since the covariates  $X_{tj}$ ,  $j = p_n + 1, \dots, p_n + d_n$ , are generated from the lags of  $Y_t$ , both temporal autocorrelation and the cross-sectional correlation among them arise. Hence, if we try to estimate or predict  $Y_t$  by running firstly the KSIS with all the potential covariates,  $X_{tj}$ ,  $j = 1, \dots, p_n + d_n$ , and secondly the PMAMAR with those that have survived the screening procedure, the results could be unsatisfactory. It is especially so when  $p_n + d_n$  is much larger than the sample size  $n$ . Due to the presence of autocorrelation in time series data, the iterative sure independence screening procedure developed in [Fan et al \(2011\)](#) cannot be applied in our context. This is because their iterative procedure involves a permutation step in which the observed data is randomly permuted to obtain a data-driven screening threshold for each iteration. When the data are autocorrelated, permutation would destroy the inherent serial dependence structure and hence may lead to erroneous thresholds being obtained. To alleviate the problem, we provide an iterative version of the KSIS+PMAMAR procedure in Appendix C of the supplementary document. This iterative procedure can be seen as a greedy selection algorithm, since at most one variable is selected in each iteration. The simulated Example 5.1 in Section 5 shows that, in general, the iterative procedure helps reduce false positive rates and increase true positive rates, especially when the exogenous covariates (i.e., the  $Z$ 's) are not correlated. This leads to the iterative procedure producing generally more accurate estimation and prediction.

## 4.2 *The PCA+KSIS+PMAMAR procedure*

In reality, the number of common factors,  $r$ , in the approximate factor model (2.11) is usually unknown. We hence need to select it from an eigenanalysis of the matrix  $\mathcal{Z}_n \mathcal{Z}_n^\top / (np_n)$ . Two ways are possible to address this issue. The first is to set a maximum number, say  $r_{\max}$  (not too large usually), for the factors. Since the factors extracted from the eigenanalysis are orthogonal to each other, the over-extracted insignificant factors will be discarded in the PMAMAR step. Another approach is to select the first few eigenvectors (corresponding to the first few largest eigenvalues) of  $\mathcal{Z}_n \mathcal{Z}_n^\top / (np_n)$  so that a pre-determined amount, say 95%, of the total variation is accounted for. See [Boneva et al \(2015\)](#) for more information on the selection of the number of common component functions. Other selection criteria such as BIC can be found in [Bai and Ng \(2002\)](#) and [Fan et al \(2013\)](#).

In the second step of the PCA+PMAMAR procedure proposed in Section 2.2, the estimated

factors and the  $d_n$  candidate lags of  $Y$  undergo a PMAMAR regression. However, since the lags of  $Y$  are often highly correlated,  $d_n$  is usually large and the PMAMAR regression usually cannot produce satisfactory results in selecting the truly significant lags. This may lead to poor performance of the PCA+PMAMAR procedure predicting the future values of  $Y$ . In order to alleviate this problem, a KSIS step can be added in between the PCA and PMAMAR steps so that the candidate lags of  $Y$  first undergo a KSIS to preliminarily screen out some insignificant lags. The simulation results in Example 5.2 below confirm that this PCA+KSIS+PMAMAR procedure improves the prediction performance of the PCA+PMAMAR procedure.

## 5 Numerical studies

In this section, we report simulation studies (Examples 5.1 and 5.2) and an empirical application (Example 5.3). Throughout this section the rule of thumb bandwidth is used as our methods do not seem to be sensitive to the choice of bandwidth.

### 5.1 Simulation studies

**Example 5.1.** In this example, the sample size is set to be  $n = 100$ , and the numbers of candidate exogenous covariates and lagged terms are  $(p_n, d_n) = (30, 10)$  and  $(p_n, d_n) = (150, 50)$ . The data-generating model is defined by

$$Y_t = m_1(Z_{t1}) + m_2(Z_{t2}) + m_3(Z_{t3}) + m_4(Z_{t4}) + m_5(Y_{t-1}) + m_6(Y_{t-2}) + m_7(Y_{t-3}) + \varepsilon_t, \quad (5.1)$$

for  $t \geq 1$ , where we set

$$\begin{aligned} m_i(x) &= -\sin(2x), \quad i = 1, 5, 6, 7, \\ m_2(x) &= x^2 - 25/12, \quad m_3(x) = x, \quad m_4(x) = e^{-x} - \frac{2}{5} \sinh(5/2), \end{aligned}$$

the exogenous covariates  $\mathbf{Z}_t = (Z_{t1}, Z_{t2}, \dots, Z_{tp_n})^\top$  are independently drawn from  $p_n$ -dimensional Gaussian distribution with zero mean and covariance matrix  $\text{cov}(\mathbf{Z}) = I_{p_n}$  or  $C_{\mathbf{Z}}$ , whose main-diagonal entries are 1 and off-diagonal entries are  $1/2$ ; the error term  $\varepsilon_t$  are independently generated from the  $\mathbf{N}(0, 0.7^2)$  distribution. The additive functions  $m_i(\cdot)$  have been chosen to be the same as those in the simulated example of Meier *et al* (2009), although they considered a static rather than a dynamic model. The real size of exogenous regressors is 4 and the real lag length is 3. We generate  $100 + n$  observations from the process (5.1) with initial states  $Y_{-2} = Y_{-1} = Y_0 = 0$  and discard the first  $100 - d_n$  observations.

The aim of this simulation is to compare the performance of the iterative KSIS+PMAMAR (IKSIS+PMAMAR) procedure in Section 4.1 with the (non-iterative) KSIS+PMAMAR procedure in Section 2.1. In order to further the comparison, we also employ the iterative sure independence

screening (ISIS) method proposed in Fan and Lv (2008), the penalized method for high-dimensional generalized additive models (penGAM) proposed in Meier *et al* (2009), and the oracle additive modelling with backfitting algorithm (Oracle, in which the true relevant variables are known). For the KSIS+PMAMAR, we choose  $\lceil 10n^{1/6} \rceil$  variables from the screening step, which then undergo a PMAMAR with the SCAD penalty. The measures of performance considered are the true positive (TP) and false positive (FP), defined, respectively, as the numbers of true and false relevant variables selected, the mean squared estimation error (MSEE) defined as  $\text{MSEE} = \frac{1}{n} \sum_{t=1}^n (Y_t - \hat{Y}_t)^2$ , where  $\hat{Y}_t$  is the fitted value of  $Y$  at  $t$  obtained from a particular method. We also generate a prediction test set of size  $n^* = 10$  and calculate 1-step-ahead forecasts for  $Y$  based on model selection and estimation from the training data set of size  $n$ . In order to compare their prediction performance, we calculate the mean squared prediction error (MSPE) of each of the methods. The MSPE is defined as  $\text{MSPE} = \frac{1}{n^*} \sum_{s=1}^{n^*} (Y_{n+s} - \hat{Y}_{n+s})^2$ , where  $\hat{Y}_{n+s}$  is the forecast of  $Y$  for time  $n + s$ . The tuning parameters in the penalized regressions are chosen by the cross-validation. The SCAD penalized regression is implemented using the R package “ncvreg”, the ISIS method implemented using the “SIS” R package, the penGAM method implemented using the “penGAM” package, and the oracle additive modelling implemented using the R package of “gam”. The results in Table 5.1 are based on 200 simulation replications, and the numbers in the parentheses are the standard errors of TP, FP, MSEE and MSPE over 200 replications.

It can be seen from Table 5.1 that when the number of candidate covariates ( $p_n + d_n = 150 + 50$ ) is much larger than the sample size, the iterative version of KSIS+PMAMAR increases the TP of the non-iterative version, and it decreases the FP in all cases except when  $\text{cov}(\mathbf{Z}) = C_{\mathbf{Z}}$  and  $(p_n, d_n) = (150, 50)$ . This results in a better performance of the IKSIS+PMAMAR than the KSIS+PMAMAR in estimation when  $(p_n, d_n) = (150, 50)$  and in prediction in all cases but when  $\text{cov}(\mathbf{Z}) = C_{\mathbf{Z}}$  and  $(p_n, d_n) = (150, 50)$ . Among the 4 variable selection procedures (i.e., IKSIS+PMAMAR, KSIS+PMAMAR, penGAM, and ISIS), the penGAM has the highest TP as well as FP. This makes it the approach that has the lowest MSEE, since within the same linear or nonlinear modelling framework it is generally the case that the more variables are selected the smaller the MSEE is. This does not hold true with MSPE. The ISIS, in contrast to the other approaches, assumes a linear modelling structure and hence is unable to correctly select truly relevant variables when the underlying data generating process is nonlinear, leading to it having the lowest TP among the 4 selection procedures. This poor performance of the ISIS in variable selection also results in its poor estimation and prediction results. The prediction performance of an approach largely depends on its accuracy in variable selection, and a low TP and high FP will lead to a high MSPE. The estimation and prediction results for the Oracle serve as a benchmark for those of the other approaches. The MSPEs from the IKSIS+PMAMAR and KSIS+PMAMAR are the closest among all the approaches to those of the Oracle. It can also be observed, by a comparison of the first two panels of Table 5.1 with the last two, that when the correlation among the exogenous variables increases, the performance of all approaches worsens.

For a fuller comparison of the above methods in this example, we have also recorded the average and median computation times for a single running of each of them. These results are presented in

Table 5.1: Average results on variable selection and accuracy of estimation and prediction in Example 5.1 over 200 replications

Model	Method	TP	FP	MSEE	MSPE
Example 5.1 $\text{cov}(\mathbf{Z}) = I_{p_n}$  $(p_n, d_n) = (30, 10)$	IKSIS+PMAMAR	6.215 (1.1983)	6.830 (4.9552)	0.9762 (0.5018)	4.1025 (2.4238)
	KSIS+PMAMAR	6.555 (0.6395)	10.350 (3.6281)	0.8559 (0.2420)	4.5385 (2.8624)
	penGAM	6.875 (0.3315)	31.030 (1.1382)	$5.6817 \times 10^{-4}$ ( $3.9801 \times 10^{-4}$ )	8.1828 (6.5490)
	ISIS	4.000 (1.0514)	8.690 (1.7113)	4.2819 (1.7572)	7.6240 (5.2296)
	Oracle	7.000 (0.0000)	0.000 (0.0000)	1.2416 (0.2192)	2.4100 (1.3857)
	IKSIS+PMAMAR	4.720 (1.4112)	10.770 (5.7727)	0.7320 (0.5010)	5.9095 (3.9042)
Example 5.1 $\text{cov}(\mathbf{Z}) = I_{p_n}$  $(p_n, d_n) = (150, 50)$	KSIS+PMAMAR	4.060 (1.2015)	13.785 (3.5781)	0.9011 (0.3200)	6.8505 (4.5468)
	penGAM	5.280 (0.8920)	56.975 (3.5180)	$8.7218 \times 10^{-5}$ ( $5.4755 \times 10^{-5}$ )	5.3088 (4.5147)
	ISIS	2.980 (0.8795)	18.020 (0.8795)	2.4720 (0.9425)	9.7411 (8.2799)
	Oracle	7.000 (0.0000)	0.000 (0.0000)	1.2709 (0.2002)	2.3877 (1.2261)
	IKSIS+PMAMAR	4.710 (1.3547)	3.605 (3.6061)	1.4600 (0.5083)	4.7438 (3.3360)
	KSIS+PMAMAR	4.950 (1.2868)	6.010 (3.9213)	1.3742 (0.4836)	5.0431 (3.5053)
Example 5.1 $\text{cov}(\mathbf{Z}) = C_{\mathbf{Z}}$  $(p_n, d_n) = (30, 10)$	penGAM	6.940 (0.2381)	31.020 (1.2992)	$9.2095 \times 10^{-4}$ ( $8.7475 \times 10^{-4}$ )	12.3246 (18.4251)
	ISIS	3.355 (1.1470)	8.280 (3.0641)	4.8761 (2.3018)	8.9032 (6.2870)
	Oracle	7.000 (0.0000)	0.000 (0.0000)	1.1927 (0.1936)	2.5930 (1.5887)
	IKSIS+PMAMAR	3.550 (1.2984)	8.600 (6.1005)	1.1042 (0.7769)	6.8042 (5.1926)
	KSIS+PMAMAR	3.115 (1.3193)	7.420 (5.5050)	1.6373 (1.0202)	6.6250 (5.6848)
	penGAM	5.575 (0.9481)	57.100 (3.6132)	$1.1801 \times 10^{-4}$ ( $1.0156 \times 10^{-4}$ )	7.6057 (6.7228)
Example 5.1 $\text{cov}(\mathbf{Z}) = C_{\mathbf{Z}}$  $(p_n, d_n) = (150, 50)$	ISIS	2.425 (0.9428)	17.820 (3.5113)	3.5018 (4.4681)	13.0225 (10.4639)
	Oracle	7.000 (0.0000)	0.000 (0.0000)	1.2428 (0.2256)	3.0629 (3.2616)

Table D.1 in Appendix D of the supplementary document, and the interested reader is referred to it for details.

**Example 5.2.** The exogenous variables  $\mathbf{Z}_t$  in this example are generated through an approximate factor model:

$$\mathbf{Z}_t = \mathbf{B}\mathbf{f}_t + \mathbf{z}_t,$$

where the rows of the  $p_n \times r$  loadings matrix  $\mathbf{B}$  and the common factors  $\mathbf{f}_t$ ,  $t = 1, \dots, n$ , are independently generated from the multivariate  $\mathbf{N}(\mathbf{0}, I_r)$  distribution, and the  $p_n$ -dimensional error terms  $\mathbf{z}_t$ ,  $t = 1, \dots, n$ , from the  $0.1\mathbf{N}(\mathbf{0}, I_{p_n})$  distribution. We set  $p_n = 30$  or  $150$ ,  $r = 3$ , and generate the response variable via

$$Y_t = m_1(f_{t1}) + m_2(f_{t2}) + m_3(f_{t3}) + m_4(Y_{t-1}) + m_5(Y_{t-2}) + m_6(Y_{t-3}) + \varepsilon_t,$$

where  $f_{ti}$  is the  $i$ -th component of  $\mathbf{f}_t$ ,  $m_1(x) = x^2 - 25/12$ ,  $m_2(x) = x$ ,  $m_3(x) = e^{-x} - \frac{2}{5} \sinh(5/2)$ ,  $m_4(x) = m_5(x) = m_6(x) = -\sin(2x)$  (these functions are the same as those in Example 5.1), and  $\varepsilon_t$ ,  $t = 1, \dots, n$ , are independently drawn from the  $\mathbf{N}(0, 0.7^2)$  distribution. In this example, we set the number of candidate lags of  $Y$  as  $d_n = 10$ . We compare the performance, in terms of estimation error and prediction error, of the following methods: PCA+PMAMAR, PCA+KSIS+PMAMAR, KSIS+PMAMAR, penGAM, ISIS, and Oracle. Since in reality both  $r$  and the factors  $\mathbf{f}_t$  are unobservable, the factors in the first two methods are estimated by the first  $\hat{r}$  eigenvectors of  $\mathcal{Z}_n \mathcal{Z}_n^\top / (np_n)$ , where  $\mathcal{Z}_n = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)^\top$ , and  $r$  is estimated by  $\hat{r}$ , where  $\hat{r}$  is chosen so that 95% of the variation in  $\mathcal{Z}_n$  is accounted for. In the PCA+PMAMAR method, the estimated factors and  $d_n$  potential lags of  $Y$  directly undergo a PMAMAR with the SCAD penalty, while in the PCA+KSIS+PMAMAR the potential lags of  $Y$  first undergo the KSIS and then the selected lags together with the estimated factors undergo a PMAMAR. The KSIS+PMAMAR, penGAM and ISIS deal directly with  $p_n$  exogenous variables in  $\mathbf{Z}_t$  and  $d_n$  lags of  $Y$  as in Example 5.1, and the Oracle uses the 3 factors and the first 3 lags, as in the true data generating process.

As in Example 5.1, the sample size is set as  $n = 100$  and the experiment is repeated for 200 times. The results are summarized in Table 5.2. It can be seen from these results that as in Example 5.1, the penGAM has the lowest MSEE but the highest MSPE as a result of it selecting a large number of variables. When the number of exogenous variables  $p_n$  is not so large compared with the sample size  $n$  (i.e., 30 compared to 100), the KSIS+PMAMAR outperforms the two PCA based approaches (i.e., PCA+PMAMAR and PCA+KSIS+PMAMAR), in terms of estimation and prediction accuracy. However, when  $p_n$  becomes larger than  $n$ , the PCA based approaches show their advantage in effective dimension reduction of the exogenous covariates, which results in their lower MSEE and MSPE. The PCA+PMAMAR has a lower MSEE but higher MSPE than the PCA+KSIS+PMAMAR. This is due to the fact that without the screening step the PCA+PMAMAR selects more false lags of  $Y$ , and the higher FP leads to a higher MSPE and lower MSEE under the same PMAMAR framework. The



Table 5.2: Accuracy of estimation and prediction in Example 5.2 over 200 replications

Model	Method	MSEE	MSPE
Example 5.2 $(p_n, d_n) = (30, 10)$	PCA+PMAMAR	1.0859 (0.3077)	5.0921 (5.5112)
	PCA+KSIS+PMAMAR	1.3553 (0.3787)	4.9089 (5.5193)
	KSIS+PMAMAR	1.0843 (0.3223)	4.3456 (6.2688)
	penGAM	0.0331 (0.0197)	21.7887 (41.1285)
	ISIS	4.7583 (2.1901)	6.8222 (6.4181)
	Oracle	1.2738 (0.2172)	2.7304 (3.3003)
Example 5.2 $(p_n, d_n) = (150, 10)$	PCA+PMAMAR	1.1465 (0.4328)	5.2313 (5.5279)
	PCA+KSIS+PMAMAR	1.4344 (0.4798)	4.9135 (4.9762)
	KSIS+PMAMAR	1.7430 (0.5134)	5.4748 (7.5642)
	penGAM	0.0053 (0.0030)	14.7217 (31.4825)
	ISIS	3.8407 (2.1800)	7.9360 (5.9894)
	Oracle	1.2753 (0.2041)	2.5206 (1.8409)

above suggests that if the focus of a study is to predict future values, there may be benefits in having the KSIS step between the PCA and PMAMAR steps to screen out insignificant lags of  $Y$ .

The computation times of the methods considered in this example are given in Table D.2 in Appendix D of the supplementary document. This table shows that the insertion of the KSIS step between PCA and PMAMAR speeds up the following PMAMAR step (as less variables undergo the PMAMAR step), leading to PCA+KSIS+PMAMAR being overall faster than PCA+PMAMAR. The interested reader is referred to Table D.2 for details.

## 5.2 An empirical application

**Example 5.3.** We next apply the proposed semiparametric model averaging methods to forecast inflation in the UK. The data were collected from the Office for National Statistics (ONS) and the Bank of England (BoE) websites and included quarterly observations on CPI and some other economics variables over the period Q1 1997 to Q4 2013. All the variables are seasonally adjusted. We use 53 series measuring aggregate real activity and other economic indicators to forecast CPI. Given the possible temporal persistence of CPI, we also add its 4 lags as predictors. Data from Q1 1997 to Q4 2012 are used as the training set and those in Q1–Q4 2013 are used for forecasting. As in [Stock and Watson \(1998, 1999\)](#), we make 4 types of transformations on different variables, depending on their nature: (i) logarithm, (ii) first difference of logarithms; (iii) first difference, and (iv) no transformation. Logarithms are usually taken on positive series that are not in rates or percentages, and first differences are taken of quantity series and of price indices. All series are standardized to have mean zero and unit variance after these transformations. Figure 5.1 plots both the original and transformed CPI series.



We use the training set to select the significant variables among the 53 exogenous economic variables and the 4 lags of CPI, as well as to estimate the model averaging weights or model coefficients. These selected variables and estimated coefficients are then used to obtain the mean squared estimation error (MSEE) and form forecasts of CPI in the four quarters of 2013. We compare the forecasting capacity of the IKSIS+PMAMAR, KSIS+PMAMAR, PCA+PMAMAR, penGAM and ISIS methods via the mean squared prediction error (MSPE) and the mean absolute prediction error (MAPE), which are defined, respectively, as

$$\text{MSPE} = \frac{1}{4} \sum_{s=1}^4 (Y_{n+s} - \widehat{Y}_{n+s})^2, \quad \text{MAPE} = \frac{1}{4} \sum_{s=1}^4 \left| Y_{n+s} - \widehat{Y}_{n+s} \right|,$$

where  $\widehat{Y}_{n+s}$  is the 1-step-ahead forecast of  $Y$  at time  $n + s$  calculated based on model selection and estimation from the training data set of size  $n$ . In addition, we also compare the estimation accuracy of the methods via the mean squared estimation error (MSEE) and the mean absolute estimation error (MAEE), defined by

$$\text{MSEE} = \frac{1}{n} \sum_{t=1}^n (Y_t - \widehat{Y}_t)^2, \quad \text{MAEE} = \frac{1}{n} \sum_{t=1}^n \left| Y_t - \widehat{Y}_t \right|,$$

where  $\widehat{Y}_t$  is the fitted value of  $Y$  at time  $t$ .

Due to the small number of candidate lags of the response ( $d = 4$ ), there is not much necessity to use the PCA+KSIS+PMAMAR approach in this example, and hence it is not included in the comparison. Similarly to [Stock and Watson \(2002\)](#), in the PCA+PMAMAR approach, common factors extracted from the exogenous variables together with lags of the response are used to forecast the response. The difference with [Stock and Watson \(2002\)](#)'s approach is that the PCA+PMAMAR allows these factors and lags to contribute to forecasting the response in a possibly nonlinear way. We also calculate forecasts based on the Phillips curve specification

$$I_{t+1} - I_t = \alpha + \beta(L)U_t + \gamma(L)\Delta I_t + \varepsilon_{t+1},$$

where  $I_t$  is the CPI in the  $t$ -th quarter,  $U_t$  is the unemployment rate,  $\beta(L) = \beta_0 + \beta_1 L + \beta_2 L^2 + \beta_3 L^3$  and  $\gamma(L) = \gamma_0 + \gamma_1 L + \gamma_2 L^2 + \gamma_3 L^3$  are lag polynomials with  $L$  being the lag operator, and  $\Delta$  is the first difference operator. We further employ some of the most commonly-used models from the BoE's suite of statistical forecasting models to model and forecast the CPI data. These include the autoregressive (AR) model, the vector autoregressive (VAR) model consisting of output, CPI, oil price, effective sterling exchange rate and BoE's base interest rate, and the smooth transition autoregressive (STAR) model. The order of autoregression in these models is selected by AIC, and the number of regimes in the STAR model is selected based on an LM test.

Table 5.3: Estimation and forecasting for UK inflation data

Method	MSEE	MSPE	MAEE	MAPE
IKSIS+PMAMAR	0.1251	0.0633	0.3019	0.2104
KSIS+PMAMAR	0.2932	0.0905	0.4476	0.2687
PCA+PMAMAR	0.1472	0.1220	0.3220	0.2434
penGAM	$1.3559 \times 10^{-5}$	0.0830	0.0029	0.2666
ISIS	0.2714	0.1037	0.4317	0.3019
Phillips Curve	1.0225	1.1900	0.7655	1.0170
AR	1.0420	0.0767	0.8011	0.2338
VAR	1.0457	0.1027	0.8287	0.2456
STAR	1.0954	0.1558	0.8361	0.2962

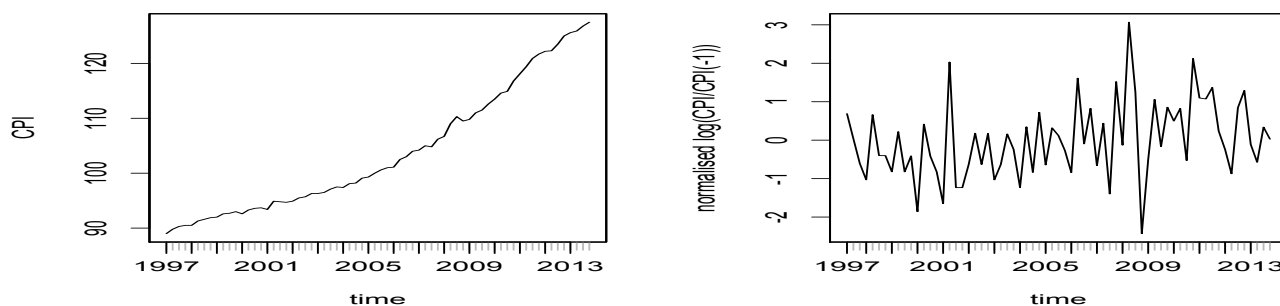


Figure 5.1: Plot of the UK CPI series. Left panel: the original UK CPI values from Q1 1997 to Q4 2013; and right panel: the normalized  $\Delta \log(\text{CPI})$ .

The MSEEs, MSPEs, MAEEs and MAPEs of the above approaches are summarized in Table 5.3, which shows that the IKSIS+PMAMAR has the smallest MSPE followed by the AR and penGAM, then KSIS+PMAMAR. The VAR and ISIS have comparable MSPEs, which are smaller than those from PCA+PMAMAR and STAR. The Phillips curve forecasts are much worse than those of the other methods. In terms of goodness of fit measured in either MSEE or MAEE, the Phillips curve, the AR, the VAR, and the STAR provide a comparable fit that is worse than that obtained from the PMAMAR based methods or the ISIS. As in the simulation studies, the penGAM gives the smallest estimation error due to a relatively large number of variables being selected. Among the variable selection/screening methods, the IKSIS+PMAMAR selects 8 exogenous variables and 2 lags of the response; the KSIS+PMAMAR selects 3 exogenous and 2 lags of response; the PCA+PMAMAR selects 14 common factors (which account for around 90% of the total variation) from the 53 exogenous variables and 3 lags of response; the penGAM selects 31 exogenous variables and 2 lags; and the ISIS selects 10 exogenous and 2 lags. Figure 5.2 provides the fitted values of the CPI observations in the training set by using the methods described above, and Figure 5.3 provides the predicted values

of the CPI from Q1 2013 to Q4 2013 using these methods. The findings from Figures 5.2 and 5.3 are consistent with those from Table 5.3. Appendix D of the supplementary document also lists the estimated models from the above methods, and the interested reader is referred to it for details.

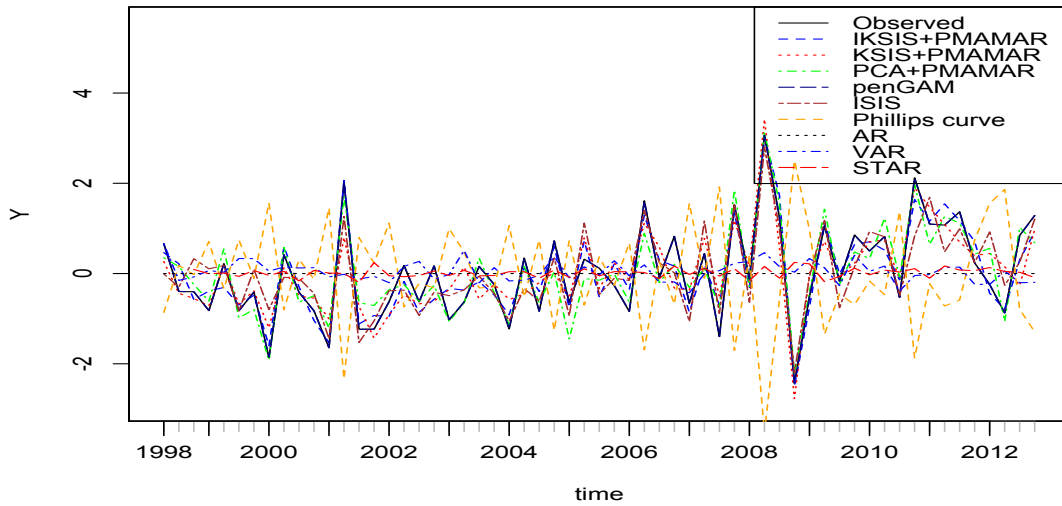


Figure 5.2: Plot of  $Y$  (normalized  $\Delta \log(\text{CPI})$ ), observed and fitted values from the methods considered.

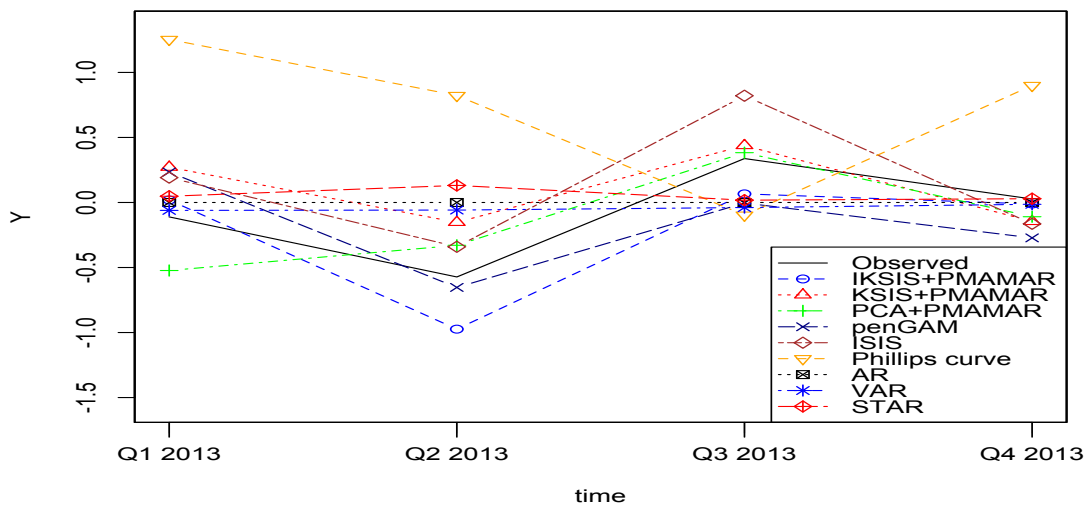


Figure 5.3: Plot of  $Y$  (normalized  $\Delta \log(\text{CPI})$ ) from Q1 2013 to Q4 2013 and their forecasts from the methods considered.

## 6 Conclusion

In this paper, we have developed two types of semiparametric methods to achieve dimension reduction on the candidate covariates and obtain good forecasting performance for the response variable. The KSIS technique, as the first step of the KSIS+PMAMAR method and the generalization of the SIS technique proposed by [Fan and Lv \(2008\)](#), screens out the regressors whose marginal regression functions do not make significant contribution to estimating the joint regression function and reduces the dimension of the regressors from an ultra large size to a moderately large size. The sure screening property developed in Theorem 1 shows that, through KSIS, the covariates whose marginal regression functions make truly significant contribution would be selected with probability approaching one. An iterative version of the KSIS is further developed in Section 4.1 and it can be seen as a possible solution to address the issue of false selection of some irrelevant covariates which are highly correlated to the significant covariates. The PMAMAR approach, as the second step of the two semiparametric dimension-reduction methods, is an extension of the MAMAR approximation introduced in [Li et al \(2015\)](#). Theorem 2 proves that the PMAMAR enjoys some well-known properties in high-dimensional variable selection such as the sparsity and oracle property. Both the simulated and empirical examples in Section 5 show that the KSIS+PMAMAR and its iterative version perform reasonably well in finite samples.

The second PCA+PMAMAR method is a generalization of the well-known factor-augmented linear regression and auto-regression models (c.f., [Stock and Watson, 2002](#); [Bernanke et al, 2005](#); [Bai and Ng, 2006](#)). By assuming an approximate factor structure on the ultra-high dimensional exogenous regressors and implementing the PCA, we estimate the unobservable factor regressors and achieve dimension reduction on the exogenous regressors. Our Theorem 3 shows that the estimated factor regressors are uniformly consistent and the asymptotic properties for the subsequent PMAMAR method (c.f., Theorem 2) remains valid for further selection of the estimated factor regressors and the time series lags. Example 5.2 shows that the PCA+PMAMAR method performs well in predicting the future value of the time series when the dimension of covariates is larger than the sample size. Furthermore, we may extend the methodology and theory developed in this paper to the more general case where some lags of the estimated factor regressors are included in the PMAMAR procedure.

## Acknowledgement

The authors are grateful to the Editor, an Associate Editor and two referees for their valuable and constructive comments which substantially improve an earlier version of the paper. Thanks also go to Dr. Lukas Meier for kindly providing the “penGAM” package and the colleagues who commented on this paper when it was presented at various conferences, workshops and research seminars. Oliver Linton’s research is supported by the Keynes fund. Zudi Lu’s research is partially supported by the Marie Curie career integration grant of European Commission.

## Supplemental document

The supplemental document contains the detailed proofs of the main asymptotic theorems given in Section 3 and some related technical lemmas. It also includes two tables recording the average and median computation times of the various methods considered in Examples 5.1 and 5.2, and lists the estimated models from the methods used in Example 5.3.

## References

- Akaike, H., 1979. A Bayesian extension of the minimum AIC procedure of autoregressive model fitting. *Biometrika* 66, 237–242.
- Ando, T., Li, K., 2014. A model averaging approach for high-dimensional regression. *Journal of the American Statistical Association* 109, 254–265.
- Bai, J., Ng, S., 2002. Determining the number of factors in approximate factor models. *Econometrica* 70, 191–221.
- Bai, J., Ng, S., 2006. Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions. *Econometrica* 74, 1135–1150.
- Bernanke, B., Boivin, J., Elias, P. S., 2005. Measuring the effects of monetary policy: a factor-augmented vector autoregressive (FAVAR) approach. *Quarterly Journal of Economics* 120, 387–422.
- Boneva, L., Linton, O., Vogt, M., 2015. A semiparametric model for heterogeneous panel data with fixed effects. *Journal of Econometrics* 188, 327–345.
- Bosq, D., 1998. *Nonparametric Statistics for Stochastic Processes: Estimation and Prediction*. Springer.
- Chamberlain, G., Rothschild, M., 1983. Arbitrage, factor structure and mean-variance analysis in large asset markets. *Econometrica* 51, 1305–1324.
- Chen, J., Li, D., Linton, O., Lu, Z., 2015. Semiparametric model averaging of ultra-high dimensional time series. *Working Paper at Institute for Fiscal Studies* available at <http://www.cemmap.ac.uk/publications/8009>.
- Chen, J., Li, D., Linton, O., Lu, Z., 2016. Semiparametric dynamic portfolio choice with multiple conditioning variables. *Journal of Econometrics* 194, 309–318.
- Cheng, M., Honda, T., Li, J., Peng, H., 2014. Nonparametric independence screening and structure identification for ultra-high dimensional longitudinal data. *Annals of Statistics* 42, 1819–1849.
- Cheng, X., Hansen, B., 2015. Forecasting with factor-augmented regression: a frequentist model averaging approach. *Journal of Econometrics* 186, 280–293.

- Claeskens, G., Hjort, N., 2008. *Model Selection and Model Averaging*. Cambridge University Press.
- Fama, E., French, K., 1992. The cross-section of expected stock returns. *Journal of Finance* 47, 427–465.
- Fan, J., Feng, Y., Song, R., 2011. Nonparametric independence screening in sparse ultra-high dimensional additive models. *Journal of the American Statistical Association* 116, 544–557.
- Fan, J., Gijbels, I., 1996. *Local Polynomial Modelling and Its Applications*. Chapman and Hall, London.
- Fan, J., Li, R., 2001. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96, 1348–1360.
- Fan, J., Liao, Y., Mincheva, M., 2013. Large covariance estimation by thresholding principal orthogonal complements (with discussions). *Journal of the Royal Statistical Society: Series B* 75, 603–680.
- Fan, J., Lv, J., 2008. Sure independence screening for ultrahigh dimensional feature space (with discussion). *Journal of the Royal Statistical Society Series B* 70, 849–911.
- Fan, J., Ma, Y., Dai, W., 2014. Nonparametric independence screening in sparse ultra-high dimensional varying coefficient models. *Journal of the American Statistical Association* 109, 1270–1284.
- Fan, J., Peng, H., 2004. Nonconcave penalized likelihood with a diverging number of parameters. *Annals of Statistics* 32, 928–961.
- Frank, I. E., Friedman, J. H., 1993. A statistical view of some chemometrics regression tools. *Technometrics* 35, 109–148.
- Green, P., Silverman, B., 1994. *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty Approach*. Chapman and Hall/CRC.
- Härdle, W., Tsybakov, A. B., 1995. Additive nonparametric regression on principal components. *Journal of Nonparametric Statistics* 5, 157–184.
- Hansen, B. E., 2007. Least squares model averaging. *Econometrica* 75, 1175–1189.
- Hansen, B. E., Racine, J., 2012. Jackknife model averaging. *Journal of Econometrics* 167, 38–46.
- Li, D., Lu, Z., Linton, O., 2012. Local linear fitting under near epoch dependence: uniform consistency with convergence rates. *Econometric Theory* 28, 935–958.
- Li, D., Linton, O., Lu, Z., 2015. A flexible semiparametric forecasting model for time series. *Journal of Econometrics* 187, 345–357.
- Liu, J., Li, R., Wu, R., 2014. Feature selection for varying coefficient models with ultrahigh dimensional covariates. *Journal of the American Statistical Association* 109, 266–274.

- Lu, Z., Linton, O., 2007. Local linear fitting under near epoch dependence. *Econometric Theory* 23, 37–70.
- Meier, L., van de Geer, S., Bühlmann, P., 2009. High-dimensional additive modeling. *Annals of Statistics* 37, 3779–3821.
- Pesaran, M. H., Pick, A., Timmermann, A., 2011. Variable selection, estimation and inference for multi-period forecasting problems. *Journal of Econometrics* 164, 173–187.
- Raftery, A. E., Madigan, D., Hoeting, J. A., 1997. Bayesian model averaging for linear regression models. *Journal of the American Statistical Association* 92, 179–191.
- Stock, J. H., Watson, M. W., 1998. Diffusion indexes. *NBER Working Paper 6702*.
- Stock, J. H., Watson, M. W., 1999. Forecasting inflation. *NBER Working Paper 7023*.
- Stock, J. H., Watson, M. W., 2002. Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association* 97, 1167–1179.
- Teräsvirta, T., Tjøstheim, D., Granger, C., 2010. *Modelling Nonlinear Economic Time Series*. Oxford University Press.
- Tibshirani, R. J., 1996. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society Series B* 58, 267–288.
- Tibshirani, R. J., 1997. The LASSO method for variable selection in the Cox model. *Statistics in Medicine* 16, 385–395.
- Wan, A. T. K., Zhang, X., Zou, G., 2010. Least squares model averaging by Mallows criterion. *Journal of Econometrics* 156, 277–283.
- Wand, M. P., Jones, M. C., 1995. *Kernel Smoothing*. Chapman and Hall.