

Time-varying Model Averaging*

Yuying Sun^{1,2,3}, Yongmiao Hong^{4,5}, Tae-Hwy Lee⁶, Shouyang Wang^{1,2,3} and Xinyu Zhang^{1,2}

¹*Academy of Mathematics and Systems Science, Chinese Academy of Sciences*

²*Center for Forecasting Science, Chinese Academy of Sciences*

³*School of Economics and Management, University of Chinese Academy of Sciences*

⁴*Department of Economics and Department of Statistics Sciences, Cornell University*

⁵*MOE Key Laboratory of Econometrics, Xiamen University*

⁶*Department of Economics, University of California, Riverside*

SUMMARY

Structural changes often occur in economics and finance due to changes in preferences, technologies, institutional arrangements, policies, crises, etc. Improving forecast accuracy of economic time series with structural changes is a long-standing problem. Model averaging aims at providing an insurance against selecting a poor forecast model. All existing model averaging approaches in the literature are designed with constant (non-time-varying) combination weights. Little attention has been paid to time-varying model averaging, which is more realistic in economics under structural changes. This paper proposes a novel model averaging estimator which selects optimal time-varying combination weights by minimizing a local jackknife criterion. It is shown that the proposed time-varying jackknife model averaging (TVJMA) estimator is asymptotically optimal in the sense of achieving the lowest possible local squared error loss in a class of time-varying model averaging estimators. Under a set of regularity assumptions, the TVJMA estimator is \sqrt{Th} -consistent. A simulation study and an empirical application highlight the merits of the proposed TVJMA estimator relative to a variety of popular estimators with constant model averaging weights and model selection.

KEY WORDS: Asymptotic optimality; Forecast combination; Local stationarity; Model averaging; Structural change; Time-varying model averaging.

JEL Classification Codes: C2, C13.

*We thank two referees, an associate editor, Oliver Linton (co-editor), Cheng Hsiao, Qingfeng Liu, Zudi Lu, Whitney Newey, Aman Ullah, Alan Tze-Kin Wan, and workshop participants at Shanghai University of Finance and Economics, UC Riverside, 2nd International Conference on Econometrics and Statistics (EcoSta 2018) at Hong Kong, 4th Guangzhou Econometrics Workshop, and 2nd Annual Forum for Chinese Econometricians at Beijing for their comments and suggestions. All remaining errors are solely ours. We acknowledge financial support from National Natural Science Foundation of China (No. 71703156) and Fujian Provincial Key Laboratory of Statistics, Xiamen University (No. 201601).

1 Introduction

Structural instability is a long-standing problem in time series econometrics (e.g., Stock & Watson (1996, 2002, 2005), Rossi (2006), and Rossi & Sekhposyan (2011)). Macroeconomic and financial time series, especially over a long period, are likely to be affected by structural instability due to changes in preferences, technologies, policies, crises, etc. For example, Stock & Watson (1996) find substantial instability in 76 representative US monthly post-war macroeconomic time series. Rossi & Sekhposyan (2011) argue that due to structural breaks, most forecast models for output growth lost their predictive ability in the mid-1970s, and became essentially useless over the last two decades. In finance, Welch & Goyal (2008) confirm that the predictive regressions of excess stock returns perform poorly in out-of-sample forecast of the U.S. equity premium, and Rapach & Zhou (2013) argue that model instability and uncertainty seriously impair the forecasting ability of individual predictive regression models. In labor economics, Hansen (2001) finds “strong evidence of a structural break in U.S. labor productivity between 1992 and 1996, and weaker evidence of a structural break in the 1960s and the early 1980s”. Thus, it is crucial to take into account such model instability and uncertainty in economic forecasting.

An approach to reducing the adverse impact of model instability and uncertainty is model averaging, which compromises across the competing models and yields an insurance against selecting a poor model. There has existed a relatively large literature on Bayesian model averaging; see Hoeting et al. (1999) for a comprehensive review. In recent years, frequentist model averaging has received growing attention in econometrics and statistics (e.g., Buckland et al. (1997), Yang (2001), Hjort & Claeskens (2003), Yuan & Yang (2005), Hansen (2007, 2008), Wan et al. (2010), Liu & Okui (2013), Liu (2015)). Most of the works focus on model averaging weights determination, related inference, and asymptotic optimality. Recently, Hansen & Racine (2012) have proposed a jackknife model averaging (JMA) which selects model averaging weights by minimizing a cross-validation criterion. The advantage of the JMA estimator mainly lies in that the asymptotic optimality theory is established under heteroskedastic error settings. Zhang et al. (2013) broaden Hansen & Racine’s (2012) scope of asymptotic optimality of the JMA estimator to encompass models with a non-spherical error covariance structure and lagged dependent variables, thus allowing for dependent data and dynamic regression models.

However, a potential problem with the aforementioned model averaging approaches is that, one predictive regression model may yield the best forecast in one period but can be dominated by other models in another period. This implies that optimal model averaging weights should change over time. There are various reasons for adopting this potentially useful time-varying approach. First, a time series model may suffer from structural instability in economics and finance. Therefore, as Stock & Watson (2003) point out, a predictor useful in one period does not guarantee its forecasting performance in other periods. The

empirical results in Stock & Watson (2007) suggest that a substantial fraction of forecasting relations are unstable. Second, macroeconomic and financial series may follow different dynamics in different time periods. For example, they may have state-dependent dynamic structures. Third, because of possible collinearity among predictors, variable selection and model selection are inherently unstable (Stock & Watson (2012)). Thus, to handle such instability, it may be better to use time-varying weights instead of constant weights in model averaging. Furthermore, since the underlying economic structure is likely to be affected by technological progress, preference changes, policy switches, crises, and so on, it is desirable to use time-varying parameter models to capture structural changes. To our knowledge, there has been no work on selecting optimal time-varying weights in model averaging where each model itself may also have time-varying parameters.

The present paper fills this gap by proposing a time-varying jackknife model averaging (TVJMA) estimator that selects model averaging weights by minimizing a local cross-validation criterion. Our approach complements the existing literature on constant JMA weights and avoids the difficulty associated with whether structural changes exist. Specifically, we assume that model parameters, as well as model averaging weights, are smooth unknown functions of time. This approach is consistent with the evidence of types of instability documented in economics, namely smooth structural changes (e.g., Rothman (1998), Grant (2002), Chen & Hong (2012) and Chen (2015)). Hansen (2001) points out that it might seem more reasonable to allow a structural change to take effect with a period of time rather than to be effective immediately. To allow the weights in model averaging to change over time, we employ the local smoothing idea to the squared error loss, leading to a local constant model averaging estimator. Moreover, we follow the spirit of Robinson (1989) and use a local constant method to estimate the time-varying parameters in each candidate model. Furthermore, we extend the candidate models from static regressions to dynamic regressions, which cover more applications in economics and finance.

In this paper, we show that the proposed TVJMA estimator is asymptotically optimal in the sense of achieving the lowest possible local squared error loss in a class of time-varying model averaging estimators, under three model settings. The first two settings admit a non-diagonal covariance structure for regression errors, including heteroscedastic errors as in Hansen & Racine (2012), with exogenous regressors. As a result, we include the non-time-varying JMA estimator in Hansen & Racine (2012) as a special case of our TVJMA estimator, under heteroscedastic errors in a nested set-up. Our theoretical analysis allows the model averaging weights to be continuously changing over time, which avoids restricting the weights to a discrete set as in Hansen & Racine (2012). The conditions required for optimality of our TVJMA estimator are neither stronger nor weaker than those required by Hansen & Racine (2012). The third model setting we consider involves lagged dependent variables with i.i.d. regression errors, where we prove the asymptotic optimality of the

TVJMA estimator by allowing the regressors to be locally stationary, in the sense of Ing & Wei (2003) and Vogt (2012).

In a simulation study and an empirical application, we compare forecast performance of the TVJMA estimator with several other model averaging estimators, including the Mallows model averaging (MMA) of Hansen (2007), JMA, a smoothed Akaike information criterion (SAIC) model averaging (Buckland et al. (1997)), a smoothed Bayesian information criterion (SBIC) model averaging, a nonparametric version of bias-corrected AIC model selection (Cai & Tiwari (2000), AICc), and a smoothed AICc (SAICc) model averaging. It is documented that for various structural changes, our TVJMA estimator outperforms these competing estimators under strictly exogenous regressors with ARMA and GARCH-type errors. Additionally, for dynamic models, the TVJMA estimator remains to be superior to other estimators under consideration.

Compared with the existing model averaging literature, our proposed approach has a number of appealing features. First, we extend conventional constant weight model averaging to time-varying weight model averaging. In particular, we propose a novel time-varying jackknife model averaging approach by exploring local information at each time point instead of over the whole sample period. The TVJMA weights selected by our method are allowed to change smoothly over time, which is consistent with evolutionary instability of economic relationships. Our result includes the constant JMA estimator in Hansen & Racine (2012) as a special case. Second, we also allow parameters in each candidate model to change smoothly over time. A nonparametric approach is used to estimate the time-varying model parameters, avoiding a potentially misspecified functional form of time-varying parameters by any parametric approach (e.g., time-varying smooth transition regression). Third, we allow regressors to be locally stationary (Dahlhaus (1996, 1997), Vogt (2012)), and as a result, time-varying parameter dynamic regression models (e.g., time-varying parameter models with lagged dependent variables) can be included as candidate models.

The remainder of this paper is organized as follows. Section 2 introduces the local jackknife criterion and develops the asymptotic optimality theory of the proposed TVJMA estimator for a general nonlinear model with heteroscedasticity. In Section 3, we consider a special class of local constant TVJMA estimators for a time-varying parameter model. Section 4 develops an asymptotic optimality theory of the TVJMA estimator for a time-varying parameter regression model with lagged dependent variables. Section 5 presents a simulation study under constant and time-varying parameter linear regressions respectively. Section 6 examines the empirical forecast performance of the TVJMA estimator for S&P 500 stock returns. Section 7 concludes. Throughout, all convergences occur when the sample size $T \rightarrow \infty$. All mathematical proofs are given in an Online Appendix.

2 Model Averaging Estimator

We consider a general nonlinear data generating process (DGP)

$$Y_t = \mu_t + \varepsilon_t = f_t(\mathbf{X}_t) + \varepsilon_t, \quad t = 1, \dots, T, \quad (1)$$

where Y_t is a dependent variable, $\mathbf{X}_t = (X_{1t}, X_{2t}, \dots)$ is possibly countably infinite, ε_t is an unobservable disturbance with $\mathbb{E}(\varepsilon_t | \mathbf{X}_t) = 0$ almost surely (a.s.), $f_t(\mathbf{x})$ is an unknown smooth function of time t , and T is the sample size. Note that when the functional form of $f_t(\cdot)$ is known up to some finite dimensional parameters, e.g., $f_t(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta}_t$, the conditional mean of Y_t given \mathbf{X}_t is parametrically specified, where parameter $\boldsymbol{\beta}_t$ is possibly time-varying. A time-varying parameter regression with $f_t(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta}_t$ will be considered in Section 3. The conventional constant parameter linear models are included as a special case if we assume that $f_t(\cdot) = f(\cdot)$ is linear. When the functional form of $f_t(\cdot)$ is unknown, we can estimate $f_t(\cdot)$ using nonparametric methods, such as the Nadaraya-Watson estimator or the local linear estimator. For notational simplicity, we let $\mathbf{Y} = (Y_1, \dots, Y_T)'$, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_T)'$ and $\mathbf{X} = (X'_1, X'_2, \dots, X'_T)'$. Furthermore, we assume that $\mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{0}$ where $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_T)'$ so that $\boldsymbol{\mu} = \mathbb{E}(\mathbf{Y} | \mathbf{X})$. We denote $\text{var}(\boldsymbol{\varepsilon} | \mathbf{X}) = \boldsymbol{\Omega}$, where $\boldsymbol{\Omega}$ is a positive definite symmetric matrix. This setup allows a non-diagonal covariance structure for regression errors. Therefore, heteroscedastic and autocorrelated errors are allowed.

2.1 Model Framework and Jackknife Criterion

Consider a sequence of candidate models indexed by $m = 1, \dots, M_T$, which are allowed to be misspecified for the underlying DGP. The number of models, M_T , may depend on the sample size T . For different models, explanatory variables may be different. Let $\{\hat{\boldsymbol{\mu}}^1, \dots, \hat{\boldsymbol{\mu}}^{M_T}\}$ be a set of nonparametric estimators of $\boldsymbol{\mu}$. Specifically, for the m -th model, the estimator of $\boldsymbol{\mu}$ may be written as $\hat{\boldsymbol{\mu}}^m = \mathbf{P}_m \mathbf{Y}$, where \mathbf{P}_m is a $T \times T$ matrix, which depends on both \mathbf{K}_t and \mathbf{X} but not on \mathbf{Y} . For instance, \mathbf{P}_m is defined in (18) below when a local constant estimator is used, and so $\hat{\boldsymbol{\mu}}^m$ is a local estimator for the conditional mean. For each time $t = 1, \dots, T$, let $\mathbf{w} = (w^1, \dots, w^{M_T})'$ be a weight vector which satisfies

$$\mathcal{H}_T = \left\{ \mathbf{w} \in [0, 1]^{M_T} : \sum_{m=1}^{M_T} w^m = 1 \right\}. \quad (2)$$

Given \mathbf{w} , an averaging estimator at any time point t for the conditional mean is

$$\hat{\mu}_t(\mathbf{w}) \equiv \sum_{m=1}^{M_T} w^m \hat{\mu}_t^m = \sum_{m=1}^{M_T} w^m \mathbf{e}_t \mathbf{P}_m \mathbf{Y} = \mathbf{e}_t \mathbf{P}(\mathbf{w}) \mathbf{Y}, \quad (3)$$

where \mathbf{e}_t is a $1 \times T$ vector, in which the t -th element is 1 and all others are zero, $\hat{\mu}_t^m = \mathbf{e}_t \mathbf{P}_m \mathbf{Y}$ and $\mathbf{P}(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \mathbf{P}_m$. Then the model averaging estimator of $\boldsymbol{\mu}$ can be fitted as $\hat{\boldsymbol{\mu}}(\mathbf{w}) = (\hat{\mu}_1(\mathbf{w}), \dots, \hat{\mu}_T(\mathbf{w}))'$.

Denote $\tilde{\boldsymbol{\mu}}^m = (\tilde{\mu}_1^m, \dots, \tilde{\mu}_T^m)'$ as the jackknife estimator of $\boldsymbol{\mu}$ for the m -th model, where $\tilde{\mu}_t^m$ is the estimator $\hat{\mu}_t^m$ obtained with the t -th observation (Y_t, \mathbf{X}_t) removed from the sample, the so-called “leave-one-out” estimator. Then, we obtain $\tilde{\boldsymbol{\mu}}^m = \tilde{\mathbf{P}}_m \mathbf{Y}$, where $\tilde{\mathbf{P}}_m$ has zeros on the diagonal and depends on \mathbf{K}_t and \mathbf{X} ; see (19) below for an example in a special setup. The jackknife model averaging estimator of μ_t , which smooths across the M_T jackknife estimators at time point t , is obtained as

$$\tilde{\mu}_t(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \tilde{\mu}_t^m = \mathbf{e}_t' \sum_{m=1}^{M_T} w^m \tilde{\mathbf{P}}_m \mathbf{Y} = \mathbf{e}_t' \tilde{\mathbf{P}}(\mathbf{w}) \mathbf{Y}, \quad (4)$$

where $\tilde{\mathbf{P}}(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \tilde{\mathbf{P}}_m$.

Set $\tilde{\boldsymbol{\mu}}(\mathbf{w}) = (\tilde{\mu}_1(\mathbf{w}), \dots, \tilde{\mu}_T(\mathbf{w}))'$. Let $\mathbf{K}_t = \text{diag}\{k_{1t}, \dots, k_{Tt}\}$, where $k_{st} = k(\frac{s-t}{Th})$, the kernel $k(\cdot) : [-1, 1] \rightarrow \mathbb{R}^+$ is a prespecified symmetric probability density, and $h \equiv h(T)$ is a bandwidth which depends on the sample size T such that $h \rightarrow 0$ and $Th \rightarrow \infty$ as $T \rightarrow \infty$. We shall minimize the local cross-validation (CV) squared error criterion,

$$\text{CV}_{t,T}(\mathbf{w}) = (\mathbf{Y} - \tilde{\boldsymbol{\mu}}(\mathbf{w}))' \mathbf{K}_t (\mathbf{Y} - \tilde{\boldsymbol{\mu}}(\mathbf{w})). \quad (5)$$

We obtain the optimal time-varying weight vector $\hat{\mathbf{w}}_t = \arg\min_{\mathbf{w} \in \mathcal{H}_T} \text{CV}_{t,T}(\mathbf{w})$, which minimizes $\text{CV}_{t,T}(\mathbf{w})$. The TVJMA estimator of μ_t for any given time point t is $\hat{\mu}_t(\hat{\mathbf{w}}_t)$.

The jackknife (or CV) criterion is widely used in selecting regression models (e.g., Allen (1974), Stone (1974) and Geisser (1975)), and the asymptotic optimality of model selection using the CV criterion is established by Li (1987) for homoskedastic regression and by Andrews (1991) for heteroskedastic regression, respectively. In this paper, the CV criterion defined above is locally weighted by \mathbf{K}_t at each time point. This local CV criterion chooses the optimal weights by generating the smallest CV value over the local sample leaving out the observation (\mathbf{X}_t, Y_t) at time t . Thus, the time-varying weight vector $\hat{\mathbf{w}}_t$ is essentially a constant weight in the neighborhood of any fixed time point t , which combines different models to yield the lowest local squared error loss.

Note that there are two key differences between our TVJMA estimator and the JMA estimators proposed by Hansen & Racine (2012) and Zhang et al. (2013). One major difference is that we allow the model averaging weights to change with time smoothly. In contrast, Hansen & Racine (2012) and Zhang et al. (2013) restrict the weights to be constant in a discrete set or a continuous set. We extend constant weights to time-varying weights, which can accommodate time-varying predictive power of candidate models. Another difference is that the models in Hansen & Racine (2012) and Zhang et al. (2013) are linear regressions, while in the present paper, $f_t(\cdot)$ can be nonlinear, and parameters in each candidate model are allowed to be unknown smooth functions of time. Theoretically, we establish the asymptotic optimality of the TVJMA estimator based on a set of smoothly time-varying parameter models. Simulation studies show that the proposed TVJMA estimator outperforms

the existing model averaging methods in the presence of smooth structural changes as well as recurrent breaks.

2.2 Asymptotic Optimality

To establish the asymptotic optimality of the TVJMA estimator, we consider the following local squared error loss and associated risk criterion:

$$L_{t,T}(\mathbf{w}) = (\hat{\boldsymbol{\mu}}(\mathbf{w}) - \boldsymbol{\mu})' \mathbf{K}_t (\hat{\boldsymbol{\mu}}(\mathbf{w}) - \boldsymbol{\mu}), \quad (6)$$

and

$$R_{t,T}(\mathbf{w}) = \mathbb{E}(L_{t,T}(\mathbf{w}) | \mathbf{X}) = \boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu} + \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega}), \quad (7)$$

where $\hat{\boldsymbol{\mu}}(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \hat{\boldsymbol{\mu}}^m$ is the weighted average of the forecasts of M_T models, and $\mathbf{A}(\mathbf{w}) = \mathbf{I}_T - \mathbf{P}(\mathbf{w})$.

Let $\tilde{L}_{t,T}(\mathbf{w})$ and $\tilde{R}_{t,T}(\mathbf{w})$ be the local jackknife squared error loss and risk, which are obtained by replacing $\hat{\boldsymbol{\mu}}(\mathbf{w})$ by $\tilde{\boldsymbol{\mu}}(\mathbf{w})$, $\mathbf{A}(\mathbf{w})$ by $\tilde{\mathbf{A}}(\mathbf{w})$, and $\mathbf{P}(\mathbf{w})$ by $\tilde{\mathbf{P}}(\mathbf{w})$, respectively. Specifically,

$$\tilde{L}_{t,T}(\mathbf{w}) = (\tilde{\boldsymbol{\mu}}(\mathbf{w}) - \boldsymbol{\mu})' \mathbf{K}_t (\tilde{\boldsymbol{\mu}}(\mathbf{w}) - \boldsymbol{\mu}) \quad (8)$$

and

$$\tilde{R}_{t,T}(\mathbf{w}) = \boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu} + \text{tr}(\tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega}). \quad (9)$$

Let

$$\xi_{t,T} = \inf_{\mathbf{w} \in \mathcal{H}_T} R_{t,T}(\mathbf{w}) \quad (10)$$

and

$$\tilde{\boldsymbol{\Omega}} = \boldsymbol{\Omega} - \text{diag}(\Omega_{11}, \dots, \Omega_{TT}), \quad (11)$$

where Ω_{tt} is the t -th diagonal element of $\boldsymbol{\Omega}$, and $\zeta(A)$ denotes the maximum singular value of matrix A .

Extending the results of Hansen & Racine (2012) and Zhang et al. (2013), we prove that the TVJMA estimator $\hat{\boldsymbol{\mu}}(\hat{\mathbf{w}}_t)$ satisfies the following optimality (OPT) property

$$(OPT) : \frac{L_{t,T}(\hat{\mathbf{w}}_t)}{\inf_{\mathbf{w} \in \mathcal{H}_T} L_{t,T}(\mathbf{w})} \xrightarrow{p} 1, \text{ as } T \rightarrow \infty.$$

This suggests that the local average squared error of the TVJMA estimator is asymptotically equivalent to the local average squared error of the infeasible best possible averaging

estimator. This optimality property is the same as that in Zhang et al. (2013), except that we now allow the weights to change smoothly over time.

To guarantee that the TVJMA estimator satisfies the OPT property under a DGP that allows smooth-changing parameters and a non-diagonal error covariance structure, we impose a set of regularity conditions:

Assumption 1. $\{\varepsilon_t\}$ is a sequence of innovations such that $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_T)'$ satisfies $\boldsymbol{\varepsilon}|\mathbf{X} \sim N(\mathbf{0}, \boldsymbol{\Omega})$, where $\boldsymbol{\Omega}$ is a $T \times T$ symmetric positive-definite matrix.

Assumption 2. The maximum singular value of $\boldsymbol{\Omega}$ satisfies $\zeta(\boldsymbol{\Omega}) \leq \bar{C} < \infty$, where \bar{C} is a constant.

Assumption 3. For $1 \leq m \leq M_T$, where M_T may depend on the sample size T , the maximum singular value of \mathbf{P}_m satisfies $\overline{\lim}_{T \rightarrow \infty} \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) < \infty$ a.s..

Assumption 4. For $1 \leq m \leq M_T$, the maximum singular value of $\tilde{\mathbf{P}}_m$ is finite when the sample size $T \rightarrow \infty$, i.e., $\overline{\lim}_{T \rightarrow \infty} \max_{1 \leq m \leq M_T} \zeta(\tilde{\mathbf{P}}_m) < \infty$ a.s..

Assumption 5. For any given time point t , the local risk $\tilde{R}_{t,T}(\mathbf{w})$, i.e., the conditional expectation of the local jackknife squared error criterion given \mathbf{X} , satisfies $\sup_{\mathbf{w} \in \mathcal{H}_T} |\tilde{R}_{t,T}(\mathbf{w})/R_{t,T}(\mathbf{w}) - 1| \rightarrow 0$ a.s. as $T \rightarrow \infty$.

Assumption 6. For any given time point t , $M_T \xi_{t,T}^{-2G} \sum_{m=1}^{M_T} R_{t,T}^G(\mathbf{w}_m^0) \rightarrow 0$ a.s., for some constant $G \geq 1$, where \mathbf{w}_m^0 is an $M_T \times 1$ weight vector with the m -th element taking the value of unity and other elements zeros.

Assumption 6'. For any given time point t , $\xi_{t,T}^{-2} \sum_{m=1}^{M_T} R_{t,T}(\mathbf{w}_m^0) \rightarrow 0$ a.s., where \mathbf{w}_m^0 is an $M_T \times 1$ weight vector with the m -th element taking the value of unity and other elements zeros.

Assumption 7. For any given time point t , $\sup_{\mathbf{w} \in \mathcal{H}_T} |\text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \tilde{\boldsymbol{\Omega}})/\tilde{R}_{t,T}(\mathbf{w})| \rightarrow 0$ a.s. as $T \rightarrow \infty$.

Assumption 8. $k : [-1, 1] \rightarrow \mathbb{R}^+$ is a symmetric bounded probability density function.

Assumption 9. The bandwidth $h = cT^{-\lambda}$ for $0 < \lambda < 1$ and $0 < c < \infty$.

Assumption 1 is the same as condition (11) in Zhang et al. (2013), which is limited to Gaussian regressions. This condition can be removed to obtain the asymptotic optimality of the TVJMA estimator for a time-varying parameter regression in Section 3. Assumption 2 ensures the largest singular values of the error covariance matrix $\boldsymbol{\Omega}$ to be finite when the sample size $T \rightarrow \infty$, corresponding to condition (12) in Zhang et al. (2013). Assumptions 3 and 4 correspond to conditions (A.3) and (A.4) of Hansen & Racine (2012), respectively.

Both of them are rather mild, because typical estimators satisfy the regularity conditions that the maximum singular values of the corresponding matrixes are bounded.

Assumption 5 imposes the condition that the leave-one-out estimator is asymptotically equivalent to the local risk of the regular estimator $\hat{\boldsymbol{\mu}}(\mathbf{w})$, uniformly over the class of averaging estimators. This is a standard condition for the application of cross-validation and almost the same as condition (10) in Zhang et al. (2013), except that the continuous time-varying set \mathcal{H}_T is used here instead of the continuous constant set \mathcal{H}_n in Zhang et al. (2013). In Section 3, we will consider time-varying parameter regressions as candidate models, where Assumption 5 is ensured by more primitive conditions; see (A.18) in Appendix.

Assumption 6 requires $M_T \sum_{m=1}^{M_T} R_{t,T}^G(\mathbf{w}_m^0) \rightarrow \infty$ at a rate slower than $\xi_{t,T}^{2G} \rightarrow \infty$ as $T \rightarrow \infty$. Assumption 6' is weaker than Assumption 6, when G is set to 1. To gain further insight into Assumptions 6 and 6', we define $\eta_{t,T} = \max_{1 \leq m \leq M_T} R_{t,T}(\mathbf{w}_m^0)$. Then, we obtain more primitive conditions for Assumptions 6 and 6' that $M_T^2 \xi_{t,T}^{-2G} \eta_{t,T}^G \rightarrow 0$ a.s. and $M_T \xi_{t,T}^{-2} \eta_{t,T} \rightarrow 0$ a.s., respectively. These conditions restrict the rates of $M_T \rightarrow \infty$, $\xi_{t,T} \rightarrow \infty$ and $\eta_{t,T} \rightarrow \infty$; in particular they require that the infimum risk $\xi_{t,T}$ explode quickly enough and the maximum risk of an individual model do not explode very quickly. Note that $\xi_{t,T} \rightarrow \infty$ is obviously necessary for Assumptions 6 and 6' to hold, which is pointed out by Hansen (2007) that this is no finite approximating model for which the bias is zero in linear regression as well as nonparametric regression. Like Ando & Li (2014), we consider a case with $\xi_{t,T} \sim T^{1-\tilde{\delta}}$ for $\tilde{\delta} < 1/2$. From Assumptions 2, 3 and 8, we can obtain $\eta_{t,T} = O_p(T)$. Given $\xi_{t,T} \rightarrow \infty$ with the rate $T^{1-\tilde{\delta}}$, $M_T \rightarrow \infty$ with a slower rate than $T^{G-\tilde{\delta}G}$ and $\eta_{t,T} = O_p(T)$, and so Assumptions 6 and 6' hold. Assumption 6 is required for the asymptotic optimality of all MMA and JMA estimators; see more discussions in Wan et al. (2010) and Zhang et al. (2013).

Assumption 7 restricts the correlation strength among unobservable disturbances and can be removed when disturbances are not correlated. Under the set-up of linear DGP, Assumption 7 can be simplified to $\sup_{\mathbf{w} \in \mathcal{H}_T} |\text{tr}(\tilde{\mathbf{P}}(\mathbf{w})\tilde{\boldsymbol{\Omega}})/\tilde{R}_{t,T}(\mathbf{w})| \rightarrow 0$ a.s. as $T \rightarrow \infty$, which is the same as condition (14) in Zhang et al. (2013). If all candidate models are linear regressions with constant parameters, it can be shown that $\sup_{\mathbf{w} \in \mathcal{H}_T} |\text{tr}(\tilde{\mathbf{P}}(\mathbf{w})\tilde{\boldsymbol{\Omega}})/\tilde{R}_{t,T}(\mathbf{w})| \leq \xi_{t,T}^{-1} \gamma \max_{1 \leq m \leq M_T} \zeta(\tilde{\mathbf{P}}_m \tilde{\boldsymbol{\Omega}})$, where γ is the number of regressors. It follows that condition (14) boils down to condition (22) in Zhang et al. (2013), which assumes that the growth rate of the number of regressors in the largest model must be slower than the rate at which $\xi_{t,T} \rightarrow \infty$. In this paper, under the linear regression setting with time-varying parameters in Section 3, we can establish the asymptotic optimality without Assumption 7.

In Assumption 8, the kernel is symmetric and bounded, and has a compact support $[-1, 1]$. It usually discounts the observations whose values are far away from the time point of interest. This implies that $k_{\max} \equiv \max_{s,t} k_{st} < \infty$, which is used in our proof. A commonly used kernel function, the Epanechnikov kernel $k(u) = 0.75(1 - u^2)I(|u| \leq 1)$, is employed in

this paper, where $I(\cdot)$ is the indicator function. Assumption 9 implies $h \rightarrow 0$ and $Th \rightarrow \infty$ as $T \rightarrow \infty$, which is a standard condition for the bandwidth; see Chen and Hong (2012). Assumption 9 includes the optimal bandwidth $h \propto T^{-1/5}$, which minimizes the integrated mean squared error (MSE) of a smoothed nonparametric estimator; see more discussions in Cai (2007) and Chen & Hong (2012).

We now state the main result of this section.

Theorem 1. *Suppose Assumptions 1-9 hold. Then for any given time point t , the TVJMA estimator $\hat{\mu}_t(\hat{\mathbf{w}}_t)$ satisfies the asymptotic optimality (OPT) property, i.e.,*

$$\frac{L_{t,T}(\hat{\mathbf{w}}_t)}{\inf_{\mathbf{w} \in \mathcal{H}_T} L_{t,T}(\mathbf{w})} \xrightarrow{p} 1.$$

Theorem 1 shows that the local squared error loss obtained from the time-varying combination weight vector $\hat{\mathbf{w}}_t$ is asymptotically equivalent to the infeasible optimal combination weight vector at any time point t . This implies that the TVJMA estimator is asymptotically optimal in the class of time-varying model averaging estimators based on possibly nonlinear models where the weight vector \mathbf{w} is restricted to the set \mathcal{H}_T , which allows the combination weights to change smoothly over time.

3 Time-varying Parameter Regression

In this section, we focus on a set of candidate models with a specific form, i.e., time-varying parameter linear regressions. This is a special case of the general candidate models in Section 2. Consider the m -th time-varying parameter regression model

$$Y_t = \mathbf{X}_t^m \boldsymbol{\beta}_t^m + \varepsilon_t^m, \quad t = 1, \dots, T, \quad m = 1, \dots, M_T, \quad (12)$$

where \mathbf{X}_t^m is a $1 \times q_m$ vector of explanatory variables, $\boldsymbol{\beta}_t^m$ is a $q_m \times 1$ possibly time-varying parameter vector, ε_t^m is an unobservable disturbance, and q_m is a positive integer that may be infinite. Note that we allow $\mathbb{E}(\varepsilon_t^m | \mathbf{X}_t^m) \neq 0$ in the set of candidate models, which arises when the m -th model is misspecified for $\mathbb{E}(Y_t | \mathbf{X}_t^m)$.

As Hansen (2001) points out, “it may seem unlikely that a structural break could be immediate and might seem more reasonable to allow a structural change to take a period of time to take effect”. We are thus interested in the following m -th smooth time-varying parameter model:

$$Y_t = \mathbf{X}_t^m \boldsymbol{\beta}^m \left(\frac{t}{T} \right) + \varepsilon_t^m, \quad t = 1, \dots, T, \quad (13)$$

where $\boldsymbol{\beta}^m : [0, 1] \rightarrow \mathbb{R}^{q_m}$ is a q_m -dimensional vector-valued function on $[0, 1]$. In the neighborhood of each time point, the model is locally stationary but it is globally nonstationary.

Various smooth time-varying parameter models have been considered to capture the evolutionary behavior of economic time series. For example, a smooth transition regression (STR) model is proposed by Chan & Tong (1986) and further studied by Lin & Teräsvirta (1994), which allows both the intercept and the slope to change smoothly over time. If the parameter function is correctly specified, parametric models for time-varying parameters can be consistently estimated with the root- T convergence rate. However, there is no economic theory to justify any concrete functional form assumption for these time-varying parameters, and the choice of a particular functional form for time-varying parameters is somewhat arbitrary, probably leading to serious misspecification. Robinson (1989, 1991) considers a nonparametric time-varying parameter model and it is further studied by Blundell et al. (1998), Cai (2007) and Chen & Hong (2012). One advantage of the nonparametric approach is that little or restrictive prior information is required for the functional forms of time-varying parameters, except for the regularity assumption that they evolve over time smoothly. In the present context, for the time-varying parameter $\beta^m(t/T)$, we follow the spirit of the smoothed nonparametric estimation in Robinson (1989).

Instead of specifying a parameterization for $\beta^m(t/T)$, which may lead to serious bias, we assume that $\beta^m(\cdot)$ is a smooth time-varying function of the ratio t/T . This assumption is based upon a common scaling scheme in the literature (e.g., Robinson (1989)). To reduce the bias and variance of a smoothed nonparametric estimator for β_t^m at any fixed time point t , it is necessary to balance the increase between the sample size T and the amount of local information at time point t . One possible solution, as suggested in Robinson (1989) and Cai (2007), is to assume a smooth function $\beta(\cdot)$ on an equally spaced grid over $[0,1]$ and consider estimation of $\beta^m(u)$ at fixed points $u \in [0,1]$. We note that the parameter β_t^m depends on the sample size T , so that new information accumulates at time point t when T increases. This ensures the consistency of parameter β_t^m at any time point t (Cai (2007), Chen & Hong (2012)).

For any s in a neighborhood of a fixed time point t , β_s^m follows a Taylor expansion:

$$\beta_s^m \approx \beta_t^m, \quad s \in [t - Th, t + Th]. \quad (14)$$

Define $\mathbf{K}_{-t} = \text{diag}\{k_{1t}, k_{2t}, \dots, k_{(t-1)t}, 0, k_{(t+1)t}, \dots, k_{Tt}\}$ as the weights for Jackknife estimation. Thus for every time point t , we obtain a local constant estimator $\hat{\beta}_t^m$ for β_t^m , and so a local least square estimator $\hat{\mu}_t^m$ and a Jackknife estimator $\tilde{\mu}_t^m$ for the m -th candidate model respectively:

$$\hat{\beta}_t^m = (\mathbf{X}^{m'} \mathbf{K}_t \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_t \mathbf{Y}, \quad (15)$$

$$\hat{\mu}_t^m = \mathbf{X}_t^m (\mathbf{X}^{m'} \mathbf{K}_t \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_t \mathbf{Y} \quad (16)$$

and

$$\tilde{\mu}_t^m = \mathbf{X}_t^m (\mathbf{X}^{m'} \mathbf{K}_{-t} \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_{-t} \mathbf{Y}. \quad (17)$$

Based on the expressions of $\hat{\mu}_t^m$ and $\tilde{\mu}_t^m$, it is straightforward to obtain

$$\mathbf{P}_m = \begin{bmatrix} \mathbf{X}_1^m (\mathbf{X}^{m'} \mathbf{K}_1 \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_1 \\ \mathbf{X}_2^m (\mathbf{X}^{m'} \mathbf{K}_2 \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_2 \\ \dots \\ \mathbf{X}_T^m (\mathbf{X}^{m'} \mathbf{K}_T \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_T \end{bmatrix} \quad (18)$$

and

$$\tilde{\mathbf{P}}_m = \begin{bmatrix} \mathbf{X}_1^m (\mathbf{X}^{m'} \mathbf{K}_{-1} \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_{-1} \\ \mathbf{X}_2^m (\mathbf{X}^{m'} \mathbf{K}_{-2} \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_{-2} \\ \dots \\ \mathbf{X}_T^m (\mathbf{X}^{m'} \mathbf{K}_{-T} \mathbf{X}^m)^{-1} \mathbf{X}^{m'} \mathbf{K}_{-T} \end{bmatrix}. \quad (19)$$

Thus, $\tilde{\mathbf{P}}_m = \mathbf{D}_m(\mathbf{P}_m - \mathbf{I}_T) + \mathbf{I}_T$, where \mathbf{D}_m is a diagonal matrix with the t -th diagonal element $(1 - h_{tt}^m)^{-1}$, and h_{tt}^m is the (t, t) element in \mathbf{P}_m .

To establish the asymptotic optimality property of $\hat{\mu}(\hat{\mathbf{w}})$, we impose the following regularity conditions:

Assumption 10. For any given time point t , $\sup_{\mathbf{w} \in \mathcal{H}_T} \text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{P}(\mathbf{w}))\xi_{t,T}^{-1} = o_p(1)$.

Assumption 11. For any given time point t , the local average of μ_t^2 is bounded, i.e., $\frac{1}{Th}\boldsymbol{\mu}'\mathbf{K}_t\boldsymbol{\mu} = O(1)$ a.s. as $T \rightarrow \infty$.

Assumption 12. For any given time point t , $h^* = O(T^{-1}h^{-1})$ and $h^{-1}\xi_{t,T}^{-1} \rightarrow 0$ a.s. as $T \rightarrow \infty$, where $\xi_{t,T}$ is defined in (10) and $h^* = \max_{1 \leq m \leq M_T} \max_{1 \leq t \leq T} h_{tt}^m$.

As pointed out by a referee, Assumption 10 implies that the bias part dominates the risk since $\text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{P}(\mathbf{w}))$ is related to the variance part of the risk. Typically, the risk is minimized by equating its bias part and its variance part. One way to make Assumption 10 hold is to restrict the set for weights. Another way is to restrict the number of candidate models or the number of variables in candidate models. Assumption 10 is the price for allowing a dependent and non-normal random error ε_t . When Assumption 1 (normal errors) is imposed or it is assumed that $\boldsymbol{\varepsilon}$ is a vector of independent variables as in the existing literature on JMA (e.g., Hansen & Racine (2012) and Ando & Li (2014)), Assumption 10 is no longer needed. Since Theorem 1 has considered normal errors, Theorem 2' below will consider the situations where $\boldsymbol{\varepsilon}$ is a vector of independent variables without using Assumption 10.

Given \mathbf{K}_t , we have $\frac{1}{Th}\boldsymbol{\mu}'\mathbf{K}_t\boldsymbol{\mu} = \frac{1}{Th}(\mu_1, \dots, \mu_T)'\mathbf{K}_t(\mu_1, \dots, \mu_T) = \frac{1}{Th} \sum_{s=1}^T k_{st}\mu_s^2 \xrightarrow{a.s.} \mathbb{E}\mu_t^2$, as $T \rightarrow \infty$. Thus, Assumption 11 implies that the local average of μ_t^2 is bounded. This is similar to condition (11) in Wan et al. (2010) and condition (23) in Zhang et al. (2013), which concern the average of μ_t^2 over the whole sample period. Finally, the first part of Assumption

12 is rather mild, which corresponds to condition (C.2) in Zhang (2015) and equation (5.2) in Andrews (1991). The second part of Assumption 12 excludes extremely unbalanced designs. This condition is reasonable and typical for the application of cross-validation; see Li (1987), Hansen & Racine (2012) and Zhang et al. (2013) for more discussions.

Theorem 2. *Suppose Assumptions 2, 3, 6' and 8-12 hold. Then for any given time point t , $\hat{\mu}_t(\hat{\mathbf{w}}_t)$ satisfies the asymptotic optimality (OPT) property.*

Theorem 2 shows that the TVJMA estimator is asymptotically optimal in the class of time-varying weighted average estimators.

Next, we establish the asymptotic optimality (OPT) result without Assumption 10. Theorem 2' below addresses the asymptotic optimality of $\hat{\mu}_t(\hat{\mathbf{w}}_t)$.

Theorem 2'. *Suppose $\boldsymbol{\varepsilon}$ is a vector of independent variables and Assumptions 2, 3, 4, 6', 8-9 and 11-12 hold. Then for any given time point t , $\hat{\mu}_t(\hat{\mathbf{w}}_t)$ satisfies the asymptotic optimality (OPT) property.*

Finally, we consider asymptotic properties of the time-varying parameter averaging estimator. Suppose the DGP is a linear time-varying parameter regression, i.e., $Y_t = \mathbf{X}_t \boldsymbol{\beta}_t + \varepsilon_t$, where \mathbf{X}_t is a $1 \times q$ vector of explanatory variables, $\boldsymbol{\beta}_t \equiv \boldsymbol{\beta}(t/T)$ is a $q \times 1$ smooth time-varying parameter vector, and $\boldsymbol{\beta} : [0, 1] \rightarrow \mathbb{R}^q$ is an unknown smooth function except for a finite number of points on $[0, 1]$. Here, q is a fixed integer, and ε_t is an unobservable disturbance with $\mathbb{E}(\varepsilon_t | \mathbf{X}_t) = 0$ almost surely. A model including only all regressors with nonzero parameters is called a true model; see Zhang (2015). Any candidate model omitting regressors with nonzero parameters is called an under-fitted model; see more discussions in Zhang (2015) and Zhang & Liu (2018). It is not required that the true model be one of the candidate models. However, at least one candidate model should not be under-fitted. This implies that one candidate model must include all these regressors with nonzero parameters and may have some redundant regressors as well. From (15), the time-varying model averaging estimator of parameter $\boldsymbol{\beta}_t$ is $\hat{\boldsymbol{\beta}}_t(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \boldsymbol{\Pi}_m' \hat{\boldsymbol{\beta}}_t^m$, where $\boldsymbol{\Pi}_m = (\mathbf{I}_{q_m}, \mathbf{0}_{q_m \times (q-q_m)})$ (i.e., a column permutation thereof) and the maximum number of columns of \mathbf{X}^m in all candidate models (i.e., $\max_{1 \leq m \leq M_T} q_m$) is bounded.

Next, we impose the following regularity conditions:

Assumption 13. *For each $j = 1, \dots, q$, the j -th element of $\boldsymbol{\beta}(\cdot)$ is continuously differentiable over the unit interval $[0, 1]$.*

Assumption 14. *For any given time point t , $\boldsymbol{\Psi}_{t,T} \equiv T^{-1} h^{-1} \sum_{s=1}^T k_{st} \mathbf{X}_s \mathbf{X}_s' \xrightarrow{p} \boldsymbol{\Psi}$ as $T \rightarrow \infty$, where $\boldsymbol{\Psi}$ is a $q \times q$ symmetric, bounded and positive definite matrix, and $T^{-1/2} h^{-1/2} \sum_{s=1}^T k_{st} \mathbf{X}_s' \varepsilon_s = O_p(1)$.*

Assumption 13 places a smoothness condition on parameters, which is commonly imposed in the literature; see Robinson (1989, 1991). Assumption 14 can be obtained from Proposition A.1 in Chen & Hong (2012) and Lemma 3 in Cai (2007). The following theorem shows that the TVJMA parameter estimator $\hat{\beta}_t(\hat{\mathbf{w}}_t)$ is \sqrt{Th} -consistent under these regularity assumptions.

Theorem 3. *Suppose Assumptions 3, 8 and 12-14 hold, and $h = cT^{-\lambda}$ for $\frac{1}{5} \leq \lambda < 1$, where $0 < c < \infty$. Then for any given time point t in the interior region $t \in [Th, T - Th]$, $\sqrt{Th}(\hat{\beta}_t(\hat{\mathbf{w}}_t) - \beta_t) = O_p(1)$ as $T \rightarrow \infty$.*

A similar result holds for the boundary regions $[1, Th] \cup [T - Th, T]$ if we assume $h = cT^{-\lambda}$ for $\frac{1}{3} \leq \lambda < 1$, where $0 < c < \infty$. This happens because the local constant estimator suffers from the well-known boundary effect problem in smoothed nonparametric estimation. As shown in Cai (2007), the convergence rate of the asymptotic bias with the local constant estimator is h^2 in the interior region, but only h in the boundary regions.

4 Asymptotic Optimality of TVJMA with Lagged Dependent Variables

In this section, we develop an asymptotic optimality theory for the TVJMA estimator based on time-varying parameter regression models that include lagged dependent variables as regressors. Dynamic regressions are widely used in macroeconomic forecasts. It is highly desirable to extend the TVJMA estimator from static regressions to dynamic regressions. Consider the following DGP

$$Y_t = \sum_{j=1}^{\infty} \beta_{jt} Y_{t-j} + \varepsilon_t, \quad t = 1, \dots, T, \quad (20)$$

where ε_t is i.i.d. with mean zero and variance σ^2 . This is a special case of the DGP in Section 2.

More generally, exogenous regressors can be added to the candidate models with finitely many lagged dependent variables. This yields an augmented regression model

$$Y_t = \sum_{j=1}^{r_1} \beta_{jt} Y_{t-j} + \sum_{j=1}^{r_2} \beta_{(r_1+j)t} X_{tj}^* + \varepsilon_t^f, \quad t = 1, \dots, T, \quad (21)$$

where X_{tj}^* is an exogenous variable, ε_t^f is the innovation, r_1 is the maximal lag order, and r_2 is the number of exogenous regressors. Let r_1 be allowed to increase and r_2 be fixed when T increases. Denote $\mathbf{Y} = (Y_1, \dots, Y_T)'$, $\mathbf{Y}_{Lt} = (Y_{t-1}, \dots, Y_{t-r_1})$, and let $\mathbf{Y}_L = (Y'_{L1}, \dots, Y'_{LT})'$ be a $T \times r_1$ matrix containing T observations of r_1 lagged dependent regressors, $\mathbf{X}^* =$

$(X_1^*, X_2^*, \dots, X_T^*)$ with $\mathbf{X}_t^* = (X_{t1}^*, \dots, X_{tr_2}^*)$ be a $T \times r_2$ matrix containing observations of r_2 exogenous regressors, $\mathbf{X} = (\mathbf{Y}_L, \mathbf{X}^*)$ be a $T \times \gamma$ matrix with rank $\gamma = r_1 + r_2$, and $\boldsymbol{\varepsilon}^f = (\varepsilon_1^f, \dots, \varepsilon_T^f)'$. The regressor matrix \mathbf{X}^m of the m -th candidate model is formed by combining the columns of \mathbf{X} . Define \mathbf{P} in a similar way to \mathbf{P}_m with \mathbf{X}^m replaced by \mathbf{X} . Note that \mathbf{X}^m is the regressor matrix in the m -th candidate model and q_m is the number of regressors in \mathbf{X}^m . Regressors are allowed to be locally stationary (Dahlhaus (1996, 1997)). Thus, our framework covers AR as well as ARX models with time-varying parameters. For each candidate model, time-varying parameters are estimated by a local constant method, which is the same as (15) in Section 3.

We impose the following regularity conditions:

Assumption 15. $\{Y_t\}$ is a locally stationary process, $\{\mathbf{X}_t^*\}$ is a strictly stationary process, and both $\{Y_t\}$ and $\{\mathbf{X}_t^*\}$ are β -mixing processes with mixing coefficients $\{\beta(j)\}$ satisfying $\sum_{j=1}^{\infty} j^2 \beta(j)^{\delta/(1+\delta)} < C < \infty$, $\sup_t \mathbb{E} \|Y_t\|^4 < C$ and $\mathbb{E} \|\mathbf{X}_t^*\|^4 < \infty$ for some constant $0 < \delta < 1$ and $C > 0$.

Assumption 16. $T q_m^{-1} h_{tt}^m = O_p(1)$, $t = 1, \dots, T$, $m = 1, \dots, M_T$, and for any given time point t , $T^{-1} h^{-1} \boldsymbol{\mu}' \mathbf{K}_t \boldsymbol{\mu} = O_p(1)$, $\gamma \xi_{t,T}^{*-1} = o_p(1)$, and $\gamma \boldsymbol{\mu}' \boldsymbol{\mu} \xi_{t,T}^{*-2} = o_p(1)$, where $\xi_{t,T}^* = \inf_{\mathbf{w} \in H_T} V_{t,T}(\mathbf{w})$ and $V_{t,T}(\mathbf{w}) = \boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu} + \sigma^2 \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}))$.

Assumption 17. For any given time point t , $\zeta(T^{-1} h^{-1} \mathbf{X}^{*'} \mathbf{K}_t \mathbf{X}^*) = O_p(1)$, $\mathbf{X}^{*'} \mathbf{K}_t \boldsymbol{\varepsilon} / \sqrt{T h} \xrightarrow{d} N(\mathbf{0}, \boldsymbol{\Delta})$, and $\zeta((T^{-1} \mathbf{X}^{*'} \mathbf{M}_t \mathbf{X}^*)^{-1}) = O_p(1)$, where $\boldsymbol{\Delta}$ is a symmetric, bounded and positive definite matrix, and $\mathbf{M}_t \equiv \mathbf{K}_t - \mathbf{K}_t \mathbf{Y}_L (\mathbf{Y}_L' \mathbf{K}_t \mathbf{Y}_L)^{-1} \mathbf{Y}_L'$.

Assumption 18. The innovation process $\{\varepsilon_t\}$ is an i.i.d. sequence with mean 0 and variance σ^2 , and satisfies that with some positive constants α_1 , α_2 and α_3 ,

$$|F_t(d_1) - F_t(d_2)| \leq \alpha_1 |d_1 - d_2|^{\alpha_2},$$

for all t when $|d_1 - d_2| \leq \alpha_3$, where $F(\cdot)$ is the distribution function of ε_t .

Assumption 19. $r_1^{6+\alpha_4} = O(T)$ for some $\alpha_4 > 0$ and $\sup_t \mathbb{E} \varepsilon_t^4 < \infty$.

In Assumption 15, local stationarity is weaker than strict stationarity. Intuitively, local stationarity implies that when the standardized time $\frac{t}{T}$ is in a neighborhood of any fixed point $\tau \in [0, 1]$, the behavior of time series $\{Y_t\}$ can be approximated up to a certain high order by a strictly stationary process $\{Y_t(\tau)\}$, and it holds that $\|Y_t - Y_t(\tau)\| = O_p(h + \frac{1}{T})$, where h is a bandwidth such that $h \rightarrow 0$ as $T \rightarrow \infty$; see Dahlhaus (1996, 1997) and Vogt (2012) for details. Thus, the autocovariance function of $\{Y_t\}$ for all times t , with $\frac{t}{T}$ in the neighborhood of τ , can be approximated arbitrarily well by that of the strictly stationary time series $\{Y_t(\tau)\}$.

Assumption 16 is analogous to Assumptions 10-12, which are used for time-varying parameter regression models when \mathbf{X}_t is assumed to be strictly stationary. The first part

of Assumption 16 is a counterpart of Assumption 12 and excludes extremely unbalanced designs. The second part of Assumption 16 concerns the local average behavior of μ_t^2 for any given time point t . Like in Shao (1997) and Wan et al. (2010), if $\{Y_t, \mathbf{X}_t\}$ is a strictly stationary process, this is the average behavior of μ_t^2 over the whole sample period. By $\boldsymbol{\mu}'\boldsymbol{\mu}/T = O_p(1)$ and Assumption 3, a sufficient condition of the fourth part of Assumption 16 is $\gamma T \xi_{t,T}^{*-2} = o_p(1)$. By comparing the expression of $V_{t,T}(\mathbf{w})$ with the risk $R_{t,T}(\mathbf{w})$ defined in (7), we can view $V_{t,T}(\mathbf{w})$ as a kind of risk as well, which may be called as a pseudo-risk. Hence, the third and fourth parts of Assumption 16 impose a restriction on the relationship among the number of regressors γ , the sample size T , and the infimum pseudo-risk $\xi_{t,T}^*$. Similar assumptions are used in Zhang et al. (2013), Liu & Okui (2013) and Ando & Li (2014).

When $\{\mathbf{X}_t' \varepsilon_t\}$ is a stationary ergodic martingale difference sequence with finite fourth moments and $T^{-1} \mathbf{X}^* \mathbf{X}'$ converges to a symmetric positive definite matrix in probability, the first part of Assumption 17 holds. The second part of Assumption 17 can be ensured by more primitive conditions; see more discussions in equation (A.7) in Cai (2007). Here, $\mathbf{X} = (\mathbf{X}_1', \dots, \mathbf{X}_T')'$, with $\mathbf{X}_t^* = (X_{t1}^*, \dots, X_{tr_2}^*)$, is a $T \times r_2$ matrix containing observations of r_2 exogenous regressors. In this paper, we assume that r_2 is fixed when T increases. It is conceivable that we could allow r_2 to increase with T at the cost of more tedious proof and other assumptions. Assumption 18 is a mild condition which is the same as condition (K.2) of Ing & Wei (2003). It holds for any distribution with a bounded probability density. This assumption is also used to prove Lemma 1 in the Mathematical Appendix. Assumption 19 is a reiteration of assumptions in Lemma 4 in the Mathematical Appendix. It can be replaced by the conditions of $r_1^{2+\alpha_4} = O(T)$ and $\sup_{-\infty < t < \infty} \mathbb{E}|\varepsilon_t|^S < \infty$ for all $S = 1, 2, \dots$.

Next, we impose conditions on the strictly stationary process $\{Y_t(\tau)\}$ indexed by $\tau \in [0, 1]$.

Assumption 20. *For any $\tau \in [0, 1]$ and $q > 0$, $\{Y_t(\tau)\}$ is strictly stationary with $E|Y_t(\tau)|^q < \infty$ and $Y_t(\tau) + \sum_{j=1}^{\infty} a_j Y_{t-j}(\tau) = \varepsilon_t, t = \dots, -1, 0, 1, \dots$, where the roots of $A(z) = 1 + \sum_{j=1}^{\infty} a_j z^j = 0$ lie outside the unit circle $|z| = 1$, and $\{\varepsilon_t\}$ is a sequence of independent random variables with mean 0 and variance σ^2 .*

Assumption 21. *For any $\tau \in [0, 1]$, $\{Y_t(\tau)\}$ is a stationary β -mixing process with mixing coefficients $\{\beta(j)\}$ satisfying $\sum_{j=1}^{\infty} j^2 \beta(j)^{\delta/(1+\delta)} < C$ for some $0 < \delta < 1$ and $0 < c < \infty$.*

Assumption 20 is a standard condition for ARMA models; see more discussions in Ing & Wei (2003). The mixing condition in Assumption 21 imposes a restriction on temporal dependence in $\{Y_t(\tau)\}$, which is commonly used in the literature (e.g., Cai (2007), Chen and Hong (2012)).

Theorem 4. *Suppose Assumptions 3, 8, 9 and 15-21 hold. Then for any given time point t , the TVJMA estimator $\hat{\mu}_t(\hat{\mathbf{w}}_t)$ in this section satisfies the asymptotic optimality (OPT) property.*

As a main contribution, Theorem 4 extends Theorem 2 for the asymptotic optimality property of the TVJMA estimator from static regression models with constant parameters to the dynamic regression models with time-varying parameters and locally stationary regressors.

5 Monte Carlo Simulation

To examine the finite sample performance of the proposed TVJMA estimator, we consider the following DGPs:

DGP 1 (Smooth Structural Changes):

$$Y_t = \mu_t + \varepsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, t = 1, \dots, T$$

where $\tau = t/T$, $F(\tau) = \tau^3$, $X_{t1} = 1$, and observations on all other regressors $\{X_{tj}, j \geq 2\}$ are generated from *i.i.d.* $N(0, 1)$ sequences. Following Hansen & Racine (2012), $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$, and the coefficient c is selected to control the population coefficient of determination $R^2 = c^2/(1 + c^2)$ to vary on a grid from 0.1 to 0.9.

To examine robustness of the TVJMA estimator, we consider three cases for $\{\varepsilon_t\}$: Case (i) $\varepsilon_t \sim i.i.d.N(0, 1)$; Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d.N(0, 1)$ and $\phi = 0.5$. This error process is the same as that of Zhang et al. (2013); Case (iii) $\varepsilon_t = \sqrt{h_t}u_t$, $h_t = 0.2 + 0.5X_{t2}^2$, $u_t \sim i.i.d.N(0, 1)$, which follows the error structure in Chen and Hong (2012). Note that $\text{var}(\varepsilon_t|X_{t2}) \neq \sigma^2$ under Case (iii).

We compare (1) the TVJMA estimator with a variety of popular model averaging estimators, namely (2) the nonparametric version of bias-corrected AIC in Cai & Tiwari (2000) (AICc); (3) a smoothed AICc (SAICc); (4) the JMA of Hansen & Racine (2012); (5) the MMA of Hansen (2007); (6) a smoothed Akaike information criterion (SAIC); and (7) a smoothed Bayesian information criterion (SBIC). The AICc for order selection is $\text{AICc} = \ln \text{RSS} + (T + \text{tr}(S^*)) / (T - (\text{tr}(S^*) + 2))$, where $\text{RSS} = \sum_{t=1}^T (Y_t - \hat{Y}_t)^2$ is based on a local constant regression and $\text{tr}(S^*)$ is the number of parameters in the model, which penalizes extra parameters for a larger value of $\text{tr}(S^*)$. For the definition of S^* , see more discussions in Cai & Tiwari (2000). The SAICc method is the model averaging estimator with the weight $w^m = \exp(-\frac{1}{2}\text{AICc}_m) / \sum_{m=1}^{M_T} \exp(-\frac{1}{2}\text{AICc}_m)$, where AICc_m is obtained from Cai & Tiwari (2000) for the m -th candidate model. The other four model averaging estimators are based on linear regressions with constant combination weights, including JMA,

MMA, SAIC and SBIC. SAIC, proposed by Buckland et al. (1997), is the least squares model averaging estimator with the weight $w^m = \exp(-\frac{1}{2}\text{AIC}_m) / \sum_{m=1}^{M_T} \exp(-\frac{1}{2}\text{AIC}_m)$, where $\text{AIC}_m = T \ln \hat{\sigma}_m^2 + 2m$. SBIC is a simplified form of the Bayesian model averaging with the weight $w^m = \exp(-\frac{1}{2}\text{BIC}_m) / \sum_{m=1}^{M_T} \exp(-\frac{1}{2}\text{BIC}_m)$, where $\text{BIC}_m = T \ln \hat{\sigma}_m^2 + m \ln T$.

The number of candidate models is determined by the rule in Hansen & Racine (2012), i.e., $M_T = \lceil 3T^{1/3} \rceil$, the nearest integer of $3T^{1/3}$. This yields $M_T = 11, 14, 15$ and 18 for $T = 50, 75, 100$ and 200 , respectively. The candidate models are $Y_t = \sum_{j=1}^m \beta_j^m(\tau) X_{tj} + \varepsilon_t^m$, $t = 1, \dots, T$, $m = 1, \dots, M_T$. For our TVJMA estimator, parameters in these candidate models are estimated by the local constant method described in Section 3. For the JMA, MMA, SAIC and SBIC methods, the parameters $\{\beta_j^m(\tau)\}$ in candidate models are assumed to be constant (i.e., they do not depend on $\tau = t/T$), and as a result, the candidate models are simplified to $Y_t = \sum_{j=1}^m \beta_j^m X_{tj} + \varepsilon_t^m$, $t = 1, \dots, T$, $m = 1, \dots, M_T$.

For the TVJMA and AICc methods, we use the Epanechnikov kernel in smoothed non-parametric estimation; this kernel has been shown to be the optimal kernel for density estimation (Epanechnikov (1969)) and robust regression (Lehmann & Casella (2006)), although our experience suggests that the choice of $k(\cdot)$ has little impact on the performance of our TVJMA estimator. For space, we report results based on a rule-of-thumb bandwidth $h = 2.34T^{-1/5}$, which attains the optimal rate for MSE (see, e.g., Chen & Hong (2012)). We generate $N = 1000$ data sets from the random sample $\{Y_t, X_t'\}_{t=1}^T$ of size T , and use the following MSE criterion to assess the accuracy of forecasts:

$$\frac{1}{N} \sum_{n=1}^N \|\hat{\boldsymbol{\mu}}(\mathbf{w})^{(n)} - \boldsymbol{\mu}^{(n)}\|^2, \quad (22)$$

where $\hat{\boldsymbol{\mu}}(\mathbf{w})^{(n)}$ and $\boldsymbol{\mu}^{(n)}$ denote the forecast value and the true value of the conditional expectation of \mathbf{Y} in the n -th replication, where $n = 1, \dots, N$. To simplify comparisons, the risk (i.e., expected squared error loss) of all model averaging estimators are normalized by the MSE of the infeasible optimal least squares model averaging estimator, which is the same as in Hansen & Racine (2012). For space, we report the Monte Calo results in graphical forms.

Figures 1-3 report the results of simulations under DGP 1. Some MSE plots are not shown in these figures, because these methods perform so poorly that their results are beyond the range of the y -axis. In most cases, the TVJMA estimator delivers the most precise forecasts among all estimators considered, especially when R^2 is relatively large. Under both conditionally heteroscedastic errors and autocorrelated errors, our method displays the best performance in terms of the risk, as is expected. Also, when the sample size T is large enough, the AICc and SAICc estimators are sometimes marginally similar to the TVJMA estimator in the cases of large R^2 . This happens because the parameters in DGP 1 are changing over time and the candidate models are time-varying parameter models as well. In most cases, the TVJMA estimator is preferred to any of the four estimators based on linear

least squares, although occasionally small to moderate reductions in MSE can be achieved for the MMA and JMA estimators with small R^2 and small T ; see $T = 50$ for example. We note that in some cases the TVJMA performances are a bit sensitive to bandwidth selection. The selection of an optimal bandwidth to estimate the time-varying combination weights is an important issue for future study. A possible solution is to consider model averaging bandwidths; see Henderson & Parmeter (2016) and Zhu et al. (2017).

Next, we consider a special case of time-varying parameter dynamic models that contain lagged dependent variables as regressors:

DGP 2 (Dynamic Regression with Smooth Structural Changes):

$$Y_t = \sum_{j=1}^{\infty} \theta_j F(\tau) Y_{t-j} + \epsilon_t,$$

where $\theta_j = 1/\sqrt{2\alpha}j^{-\alpha-1/2}$, $F(\tau) = \tau$, $\epsilon_t = \frac{\sqrt{R^2}}{c}\varepsilon_t$, $c = \sum_{j=1}^{\infty} \theta_j^2$, $\varepsilon_t \sim i.i.d.N(0, 1)$ and $\alpha = 1.5$. We allow R^2 to vary on a grid from 0.1 to 0.9.

Furthermore, to investigate the finite sample performance of the TVJMA estimator under DGPs with various structural changes, we consider following three DGPs with Case (ii) for $\{\varepsilon_t\}$. For DGPs 3-5 below, $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with various values of $c > 0$ and $\alpha = 1.5$. These parameter values are the same as those in DGP 1:

DGP 3 (Single Structural Break):

$$Y_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t,$$

where $F(\tau) = 0.5I(\tau \leq 0.3) + I(\tau > 0.3)$ and $\tau = t/T$.

DGP 4 (Smooth Transition Regression):

$$Y_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t,$$

where $F(\tau) = 1.5 - 1.5 \exp(-3(\tau - 0.3)^2)$ and $\tau = t/T$.

DGP 5 (Smooth Structural Changes with Periodicity):

$$Y_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t,$$

where $F(\tau) = \sin(\pi\tau^2)$ and $\tau = t/T$.

For each of DGPs 2-5, we generate N data sets of the random sample $\{X_t, Y_t\}_{t=1}^T$ for each sample size $T = 50, 75, 100$ and 200 , where $X_{t1} = 1$ and observations on all other regressors $\{X_{tj}, j \geq 2\}$ are generated from $i.i.d.N(0, 1)$ sequences. The candidate models and their parameter estimation methods under DGPs 2-5 are the same as those under DGP 1. Specifically, DGP 2 is a dynamic linear regression model with time-varying parameters,

which is based on Section 4. DGPs 3-5 are based on the same set-up as that of DGP 1, except that DGPs 3-5 focus on various structural changes with Case (ii) for $\{\varepsilon_t\}$. The results are reported in Figures 4-7.

In Figure 4, we consider the dynamic regression model with smooth time-varying parameters under DGP 2. When the sample size T is large enough, the TVJMA estimator yields a smaller risk than all other four estimators. This is even more clear for small R^2 .

In Figure 5, we consider the deterministic single break under DGP 3, namely, a single break with a given breakpoint and size. The TVJMA estimator, not surprisingly, outperforms all other estimators when the sample size T is larger than 50 for all R^2 , while AICc and SAICc yield smaller risks than SAIC and SBIC respectively; see, for example, the case with $T = 200$ and $R^2 > 0.4$.

In Figure 6, we consider the smooth transition regression with nonmonotonic smooth structural changes under DGP 4. This is considered in Lin & Teräsvirta (1994), which is further studied by Cai (2007) and Chen (2015). The smooth transition function is a second-order logistic function. The TVJMA estimator dominates all other estimators. We note that in most cases, the AICc estimator is similar to the SAICc estimator for large T and large R^2 , while both of them have a higher risk than the TVJMA estimator. SAIC achieves a lower risk for a smaller R^2 and SBIC is the least accurate estimator for large R^2 .

In Figure 7, we consider DGP 5, which has periodic structural changes, covering long or short period cycles; see Twrđy & Batista (2016) for an example of container throughput forecasting. The TVJMA estimator outperforms all other estimators. The SAICc estimator is the worst performing estimator when $R^2 < 0.3$, while its performance improves as R^2 increases and yields the second smallest risk when $R^2 \geq 0.7$.

To sum up, the TVJMA estimator achieves the lowest risk among all the model averaging estimators under various DGPs. When the sample size T increases, even for small R^2 , the TVJMA appears to be the best estimator. When R^2 is large, the SAICc estimator achieves a lower risk than the AICc model selection, which is consistent with the findings in the earlier literature. However, both of them perform worse than the TVJMA estimator for large T and all R^2 . We also consider a benchmark nonparametric local constant estimator without any model selection. It is shown that the local constant model without model selection performs quite poorly relative to other methods in most cases. Furthermore, following a referee's suggestion, we also compare the TVJMA estimator with a time-varying leave- k -out cross-validation model averaging (LkoMA) method (e.g., Gao et al. (2016)). We find that when R^2 is small, the TVJMA estimator outperforms the time-varying LkoMA estimator under different DGPs, especially DGP 2. Nevertheless, when R^2 is large, the time-varying LkoMA estimator achieves a slightly lower risk than the TVJMA estimator except for DGP 2. Developing optimal time-varying leave- k -out cross-validation weight selection methods and extending the proof technique for the asymptotic optimality property are important

topics for future research.

6 Empirical Application

It is widely accepted that stock return predictability is an important yet controversial issue in empirical finance. The conventional wisdom, studied by Campbell (1990) and Cochrane (1996), is that aggregate dividend yields strongly forecast excess stock return, even at longer horizons. Other commonly used predictive variables are financial ratios, such as dividend-price ratio, earnings-price ratio, and book-to-market ratio (Rozeff (1984), Fama & French (1988), Campbell & Shiller (1988), Lewellen (2004)), as well as corporate payout and financing activity (Lamont (1998), Baker & Wurgler (2000)). However, Wang (2003) and Welch & Goyal (2008) show that predictive regressions of excess stock returns perform poorly in out-of-sample forecasts of the U.S. equity premium while historical average returns generate superior forecasts, which causes vigorous debates in the literature (Campbell & Thompson (2008)). It is possible that the presence of structural changes leads to a changing predictive relationship. Indeed, Pesaran & Timmermann (2007) find that the size of parameter variations between the break points in models is considerably large, and the parameter estimates of dividend yields take even opposite signs before and after 1991. Chen & Hong (2012) find strong evidence against stability in univariate and multivariate predictor regressions for both the postwar and post-oil-shock sample periods. Furthermore, Rapach & Zhou (2013) point out that model instability and uncertainty seriously impair the forecasting ability of predictive regression models.

The sensitivity of empirical results to model parameter estimation highlights the need of time-varying combination weights in model averaging. In this section, we compare the performance of stock return forecasts using our TVJMA method and existing methods. The key distinction between these methods lies in that we allow model combination weights to change over time in combining time-varying parameter predictive models.

We employ Campbell and Thompson’s (2008) popular dataset, which is used in Chen & Hong (2012), Jin et al. (2014) and Lu & Su (2015), among many others. We consider the following predictive regression model:

$$Y_{t+1} = \alpha_t + \beta_t' \mathbf{X}_t + \varepsilon_{t+1},$$

where $Y_{t+1} = \ln[(P_{t+1} + D_{t+1})/P_t] - r_t$, P_t is the S&P 500 price index, D_t is the dividend paid on the S&P 500 price index, r_t is the 3-month treasury bill rates, \mathbf{X}_t is a set of predictive variables, i.e., $\mathbf{X}_t = (X_{t1}, \dots, X_{tp})'$, and p is the number of predictive variables. Quarterly variables from Welch & Goyal (2008) are available for 1927:01-2005:12, since quarterly stock returns before 1927 are constructed by interpolation of lower-frequency data, which may be not reliable.

Following Welch & Goyal (2008) and Rapach et al. (2010), we consider 14 financial and economic variables, sorted by relevance to \mathbf{Y} : default yield spread (X_1), treasury bill rate (X_2), net equity expansion (X_3), term spread (X_4), log dividend price ratio (X_5), log earnings price ratio (X_6), long-term yield (X_7), book-to-market ratio (X_8), inflation (X_9), log dividend yield (X_{10}), log dividend payout ratio (X_{11}), stock variance (X_{12}), long-term return (X_{13}), default return spread (X_{14}). For simplicity, we consider the following 14 nested candidate models: $\{1, X_1\}$, $\{1, X_1, X_2\}$, \dots , $\{1, X_1, \dots, X_{14}\}$. All candidate models are time-varying parameter linear regression models, and parameters are estimated by the local constant method in (15) in Section 3.

The estimation sample starts from 1947Q1 and our estimation is based on subsamples with size $T_1 = 80, 92, 104, 116, 128, 140, 152, 164, 176, 188, 200, 212$ and 224 , respectively. The remaining observations are used for out-of-sample recursive forecast accuracy assessment. For example, we use the model averaging weights for the time period T_1 , $\hat{\mathbf{w}}_{T_1}$, to construct a forecast of Y_{T_1+1} . After that we input a new observation and recalculate new model averaging weights for the time period $T_1 + 1$ and then obtain a forecast of Y_{T_1+2} . Thus, the out-of-sample forecast periods begins from 1967Q1, 1970Q1, 1973Q1, 1976Q1, 1979Q1, 1982Q1, 1985Q1, 1988Q1, 1991Q1, 1994Q1, 1997Q1, 2000Q1 and 2003Q1, respectively, and all end at 2005Q4. The postwar sample, covering 1947Q1-2005Q4, and the post-oil-shock subsample, covering 1976Q1-2005Q4, are commonly used in the literature, e.g., Welch & Goyal (2008), Chen & Hong (2012), etc. The bandwidth employed in TVJMA, AICc and smoothed AICc is set to be $2.34T_1^{-0.2}$. Following Ullah et al. (2017), we use the out-of-sample \tilde{R}^2 measure:

$$\tilde{R}^2 = 1 - \frac{\sum_{t=T_1}^{T-1} (Y_{t+1} - \hat{Y}_{t+1})^2}{\sum_{t=T_1}^{T-1} (Y_{t+1} - \bar{Y})^2},$$

where \hat{Y}_{t+1} is the prediction of Y_{t+1} based on a given forecast method, and \bar{Y} is the historical average of Y_t over the T_1 observations. This measure represents the relative difference in squared error predictive risks. The negative (positive) value of \tilde{R}^2 suggests that \hat{Y} yields a larger (smaller) sum of squared one-period forecast errors than the historical average method.

Another measure we use is the mean square predictive error (MSPE), which is widely used in the literature (e.g., Sun et al. (2018)).

$$\text{MSPE} = \frac{1}{T - T_1} \sum_{t=T_1}^{T-1} (Y_{t+1} - \hat{Y}_{t+1})^2. \quad (23)$$

Tables 1 and 2 compare \tilde{R}^2 and MSPE between the TVJMA estimator and other estimators. We find that in most cases, the TVJMA estimator is almost always the best estimator among all methods considered. Our finding supports the argument of Chen & Hong (2012) that instability exists in univariate predictor models for stock returns and smooth structural

change is a possibility, which explains why the TVJMA estimator is more appropriate than JMA and MMA. We note that the JMA estimator yields the second smallest forecast errors in most cases, with the MMA estimator being a close fourth. In most cases, the AICc estimator yields the worst performance. It is possible that the evidence of instability is a bit weak in quarterly data, which is consistent with the findings in Chen & Hong (2012).

7 Conclusion

Although structural changes have received considerable attention in time series econometrics for a long time, no work has attempted to consider time-varying model averaging for both linear and nonlinear candidate models, including those with time-varying parameters. We propose a frequentist method for model averaging with time-varying jackknife combination weights. This method is more appropriate than the conventional MMA and JMA methods under structural changes. It is shown that our TVJMA estimator is asymptotically optimal in the sense of achieving the lowest possible squared error loss in a class of time-varying model average estimators. In a simulation study, we document that the TVJMA method outperforms a variety of existing methods, including a nonparametric version of the bias-correct AIC method. An application to predicting stock returns also demonstrates that the TVJMA method outperforms many model averaging methods.

We conclude this paper by pointing out some important areas of future work. First, it would be interesting to propose a time-varying lasso-type method to select relevant regressors from a set of many potential predictive variables in the first step, and then consider time-varying model averaging in the second step. This would allow different sets of regressors (so different models) in different time periods. In these scenarios, time-varying model averaging weights are expected to yield robust and accurate forecasts. Second, this paper has only considered a global bandwidth for the TVJMA estimator, which may be severely affected by the existence of structural changes. It will be desirable to use a time-varying bandwidth for each time point. Finally, an extension of “leave- k -out” cross-validation model averaging (e.g., Gao et al. (2016)) to allow for time-varying combination weights would be highly interesting, which may be more appropriate for averaging time series models under structural changes.

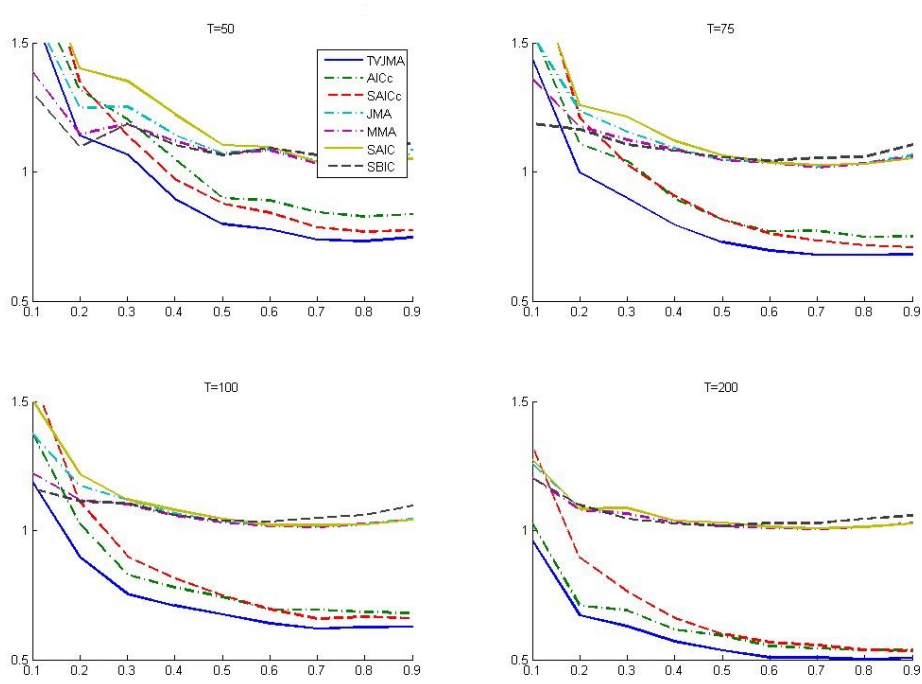


Figure 1: Finite-sample Performance under DGP 1 with Case (i)

Notes: (1) DGP 1 (Smooth Structural Changes):

$$Y_t = \mu_t + \varepsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, \quad t = 1, \dots, T,$$

where $\tau = t/T$, $F(\tau) = \tau^3$, $X_{t1} = 1$, and all other regressors $\{X_{tj}, j \geq 2\}$ are *i.i.d.* $N(0, 1)$ sequences; $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$.

(2) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(3) Three cases for $\{\varepsilon_t\}$: Case (i) $\varepsilon_t \sim i.i.d.N(0, 1)$; Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d.N(0, 1)$ and $\phi = 0.5$; Case (iii) $\varepsilon_t = \sqrt{h_t}u_t$, $h_t = 0.2 + 0.5X_{t2}^2$, $u_t \sim i.i.d.N(0, 1)$.

(4) In each panel, the y -axis and the x -axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

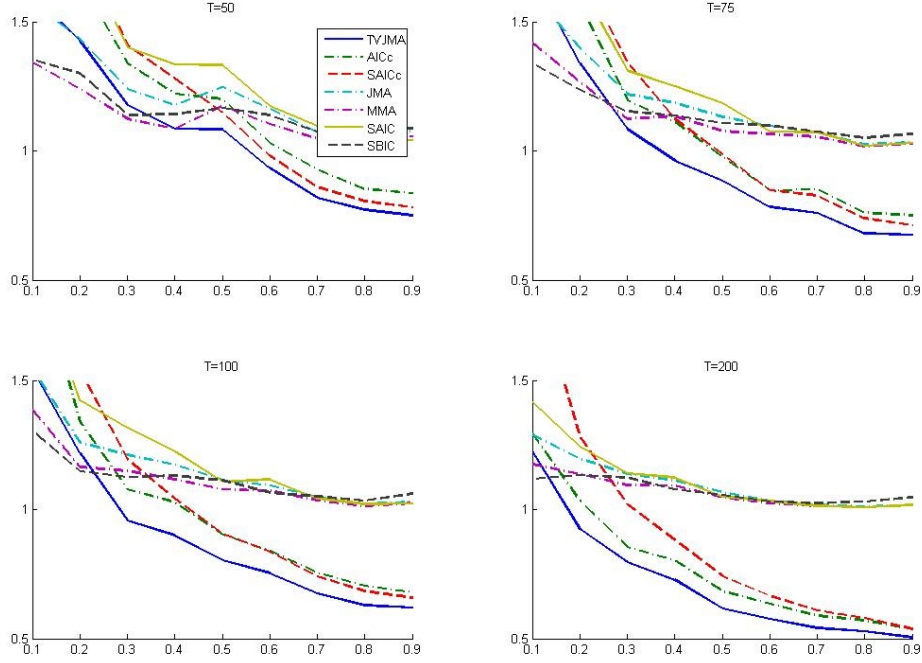


Figure 2: Finite-sample Performance under DGP 1 with Case (ii)

Notes: (1) DGP 1 (Smooth Structural Changes):

$$Y_t = \mu_t + \varepsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, \quad t = 1, \dots, T,$$

where $\tau = t/T$, $F(\tau) = \tau^3$, $X_{t1} = 1$, and all other regressors $\{X_{tj}, j \geq 2\}$ are *i.i.d.* $N(0, 1)$ sequences; $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$.

(2) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(3) Three cases for $\{\varepsilon_t\}$: Case (i) $\varepsilon_t \sim i.i.d.N(0, 1)$; Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d.N(0, 1)$ and $\phi = 0.5$; Case (iii) $\varepsilon_t = \sqrt{h_t}u_t$, $h_t = 0.2 + 0.5X_{t2}^2$, $u_t \sim i.i.d.N(0, 1)$.

(4) In each panel, the y axis and the x axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

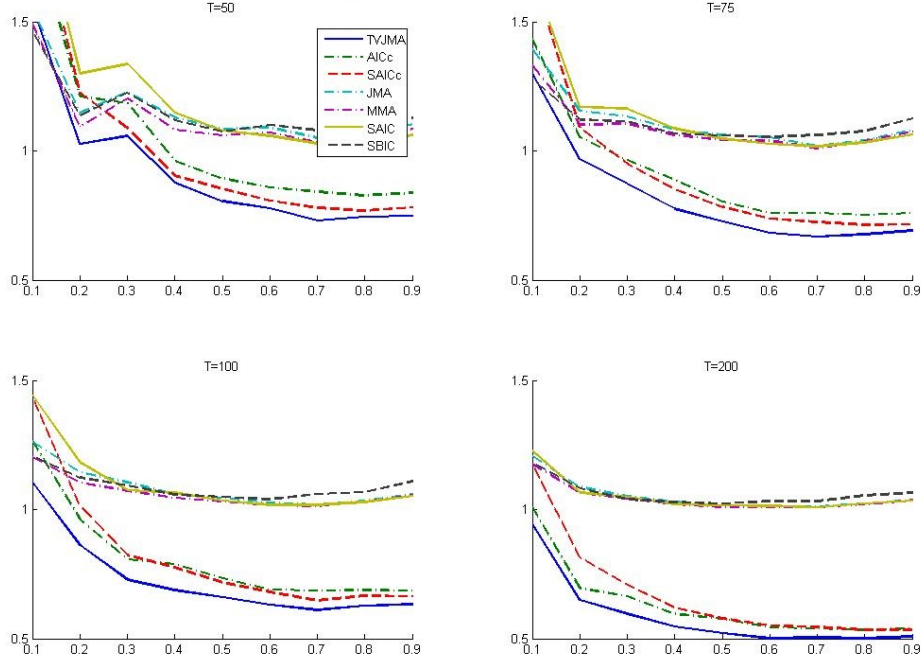


Figure 3: Finite-sample Performance under DGP 1 with Case (iii)

Notes: (1) DGP 1 (Smooth Structural Changes):

$$Y_t = \mu_t + \varepsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, \quad t = 1, \dots, T,$$

where $\tau = t/T$, $F(\tau) = \tau^3$, $X_{t1} = 1$, and all other regressors $\{X_{tj}, j \geq 2\}$ are *i.i.d.* $N(0, 1)$ sequences; $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$.

(2) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(3) Three cases for $\{\varepsilon_t\}$: Case (i) $\varepsilon_t \sim i.i.d.N(0, 1)$; Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d.N(0, 1)$ and $\phi = 0.5$; Case (iii) $\varepsilon_t = \sqrt{h_t}u_t$, $h_t = 0.2 + 0.5X_{t2}^2$, $u_t \sim i.i.d.N(0, 1)$.

(4) In each panel, the y-axis and the x-axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

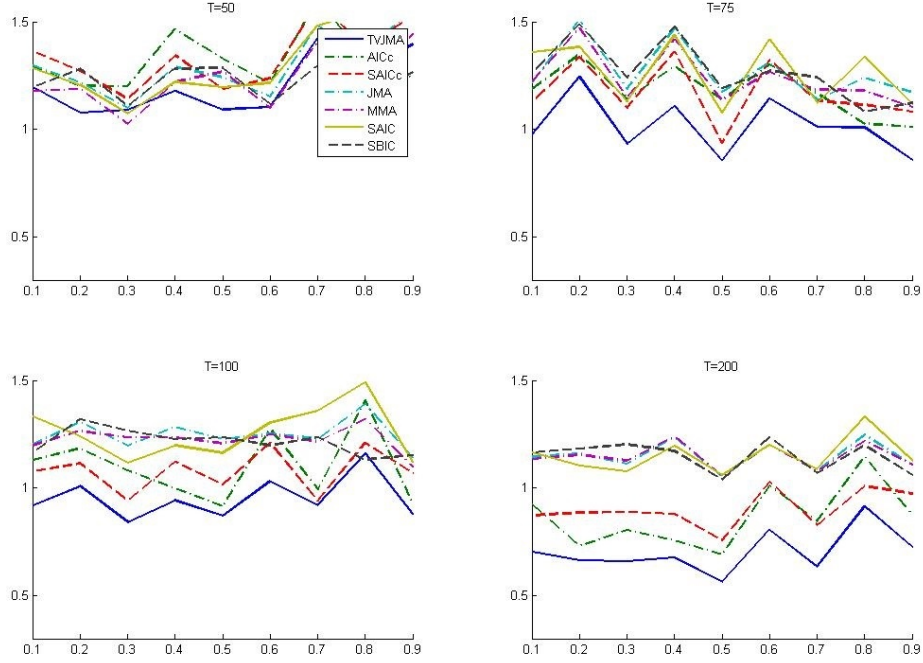


Figure 4: Finite-sample Performance under DGP 2

Notes: (1) DGP 2 (Dynamic Regression with Smooth Structural Changes):

$$Y_t = \mu_t + \epsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) Y_{t-j} + \epsilon_t, \quad t = 1, \dots, T,$$

where $\theta_j = 1/\sqrt{2\alpha}j^{-\alpha-1/2}$, $c = \sum \theta_j^2$, $F(\tau) = \tau$, $\epsilon_t = \frac{R}{c}\varepsilon_t$ with R^2 varying on a grid from 0.1 to 0.9, $\varepsilon_t \sim i.i.d.N(0, 1)$ and $\alpha = 1.5$.

(2) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(3) In each panel, the y -axis and the x -axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

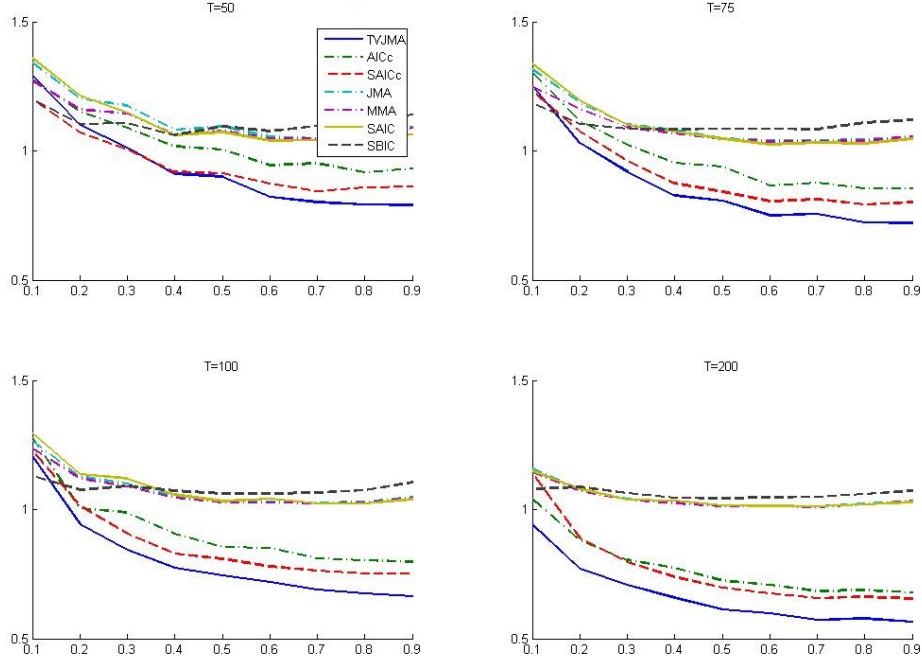


Figure 5: Finite-sample Performance under DGP 3 with Case (ii)

Notes: (1) DGP 3 (Single Structural Break):

$$Y_t = \mu_t + \varepsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, \quad t = 1, \dots, T,$$

where $F(\tau) = 0.5I(\tau \leq 0.3) + I(\tau > 0.3)$, $\tau = t/T$, $X_{t1} = 1$, and all other regressors $\{X_{tj}, j \geq 2\}$ are *i.i.d.* $N(0, 1)$ sequences; $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$.

(2) In DGP 3, $\{\varepsilon_t\}$ follows Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d. N(0, 1)$ and $\phi = 0.5$.

(3) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(4) In each panel, the y -axis and the x -axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

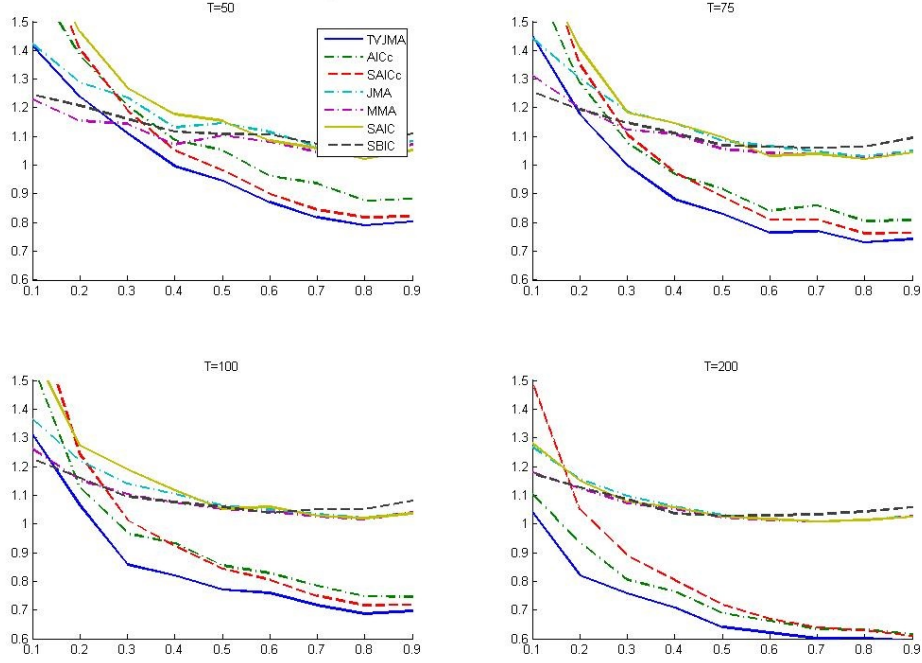


Figure 6: Finite-sample Performance under DGP 4 with Case (ii)

Notes: (1) DGP 4 (Smooth transition regression):

$$Y_t = \mu_t + \varepsilon_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, \quad t = 1, \dots, T,$$

where $F(\tau) = 1.5 - 1.5 \exp(-3(\tau - 0.3)^2)$, $\tau = t/T$, $X_{t1} = 1$ all other regressors $\{X_{tj}, j \geq 2\}$ are *i.i.d.* $N(0, 1)$ sequences; $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$.

(2) In DGP 4, $\{\varepsilon_t\}$ follows Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d. N(0, 1)$ and $\phi = 0.5$.

(3) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(4) In each panel, the y -axis and the x -axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

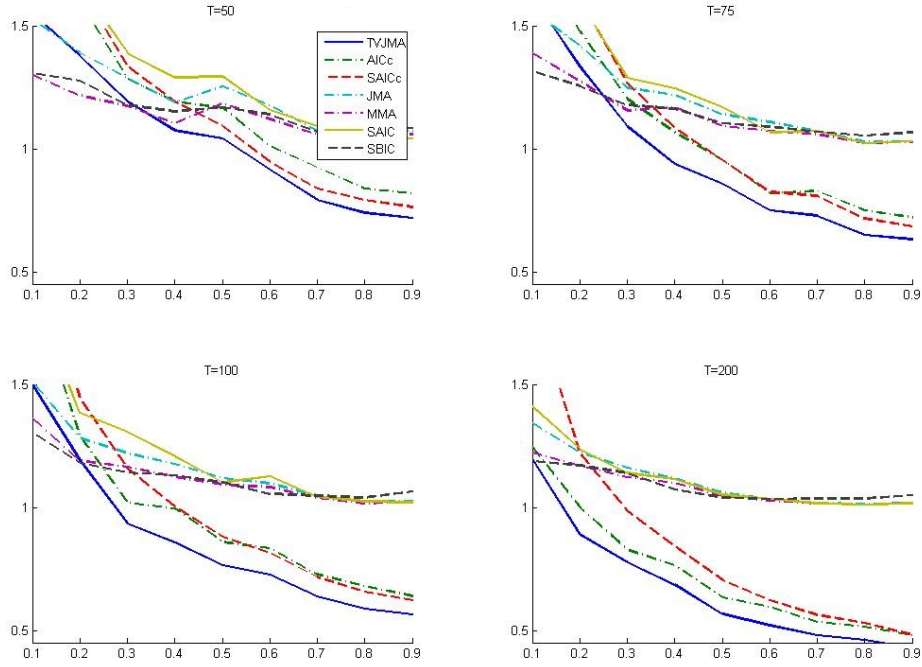


Figure 7: Finite-sample Performance under DGP 5 with Case (ii)

Notes: (1) DGP 5 (Smooth Structural Changes with Periodicity):

$$Y_t = \sum_{j=1}^{\infty} \theta_j F(\tau) X_{tj} + \varepsilon_t, \quad t = 1, \dots, T,$$

where $F(\tau) = \sin(\pi\tau^2)$, $\tau = t/T$, $X_{t1} = 1$, and all other regressors $\{X_{tj}, j \geq 2\}$ are *i.i.d.* $N(0, 1)$ sequences; $\theta_j = c\sqrt{2\alpha}j^{-\alpha-1/2}$, with $c > 0$ and $\alpha = 1.5$.

(2) In DGP 5, $\{\varepsilon_t\}$ follows Case (ii) $\varepsilon_t = e_{t,1} + e_{t,2}$, $e_{t,1} \sim N(0, X_{t2}^2)$, $e_{t,2} = \phi e_{t-1,2} + u_t$, $u_t \sim i.i.d.N(0, 1)$ and $\phi = 0.5$.

(3) In each figure, the sample sizes are shown in four panels. The sample size varies from $T = 50, 75, 100$ and 200 .

(4) In each panel, the y -axis and the x -axis display the MSE and the population R^2 , respectively. Seven methods to estimate parameters are shown in these figures, including TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators.

Table 1: Out-of-sample \tilde{R}^2 of Different Methods

		TVJMA	AICc	SAICc	JMA	MMA	SAIC	SBIC
Estimation	Prediction	$h = 2.34T_1^{-0.2}$						
1947Q1($T_1=80$)	1967Q1	0.1771°	0.0335	0.1018	0.1761°°	0.1657	0.1111	0.1024
1947Q1($T_1=92$)	1970Q1	0.1242°	-0.0312	0.0549	0.1228°°	0.1124	0.0512	0.0530
1947Q1($T_1=104$)	1973Q1	0.1212°	-0.0443	0.0770	0.1107°°	0.0977	0.0358	0.0367
1947Q1($T_1=116$)	1976Q1	0.0372°	-0.1574	-0.0397	0.0025°°	-0.0165	-0.1128	-0.0708
1947Q1($T_1=128$)	1979Q1	0.0357°	-0.1727	-0.0291	-0.0188°°	-0.0381	-0.1393	-0.1007
1947Q1($T_1=140$)	1982Q1	-0.1057°	-0.3579	-0.1375°°	-0.1830	-0.2064	-0.3210	-0.2220
1947Q1($T_1=152$)	1985Q1	-0.1833°	-0.4899	-0.2359°°	-0.2586	-0.2829	-0.4043	-0.2567
1947Q1($T_1=164$)	1988Q1	-0.2630°°	-0.6292	-0.2522 °	-0.3773	-0.4091	-0.5724	-0.3658
1947Q1($T_1=176$)	1991Q1	-0.2238°°	-0.4310	-0.2242	-0.2194°	-0.2347	-0.2945	-0.3095
1947Q1($T_1=188$)	1994Q1	-0.2181	-0.4080	-0.1579°	-0.2207	-0.2161°°	-0.2570	-0.3249
1947Q1($T_1=200$)	1997Q1	-0.0423°°	-0.1979	-0.1005	-0.0365°	-0.0538	-0.0735	-0.0990
1947Q1($T_1=212$)	2000Q1	0.0125°°	-0.1434	-0.1316	0.0694°	-0.0207	-0.0383	-0.0330
1947Q1($T_1=224$)	2003Q1	0.1859	-0.0714	-0.0560	0.1822	0.2909°°	0.3013°	-0.0543

Notes: (1) The estimation sample begins from 1947Q1, with T_1 observations. The prediction period begins from the quarter indicated in the second column.

(2) Seven methods are shown in Table 1: TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators. The larger the criteria, the better the method.

(3) The bandwidth used here is $2.34T_1^{-0.2}$, the same as that in the simulation study.

(4) ° and °° denote the best and the second best forecast among these seven methods, respectively.

Table 2: Out-of-sample MSPE of Different Methods

Estimation	Prediction	TVJMA	AICc	SAICc	JMA	MMA	SAIC	SBIC
		$h = 2.34T_1^{-0.2}$						
1947Q1($T_1=80$)	1967Q1	0.0758°	0.0891	0.0827	0.0759°°	0.0769	0.0819	0.0827
1947Q1($T_1=92$)	1970Q1	0.0803°	0.0945	0.0866	0.0804°°	0.0814	0.0870	0.0868
1947Q1($T_1=104$)	1973Q1	0.0797°	0.0947	0.0837	0.0806°°	0.0818	0.0874	0.0874
1947Q1($T_1=116$)	1976Q1	0.0693°	0.0833	0.0748	0.0717°°	0.0731	0.0800	0.0770
1947Q1($T_1=128$)	1979Q1	0.0708°	0.0861	0.0755	0.0748°°	0.0762	0.0836	0.0808
1947Q1($T_1=140$)	1982Q1	0.0740°	0.0908	0.0761°°	0.0791	0.0807	0.0884	0.0817
1947Q1($T_1=152$)	1985Q1	0.0779°	0.0981	0.0814°°	0.0829	0.0845	0.0925	0.0828
1947Q1($T_1=164$)	1988Q1	0.0695°°	0.0897	0.0690°	0.0758	0.0776	0.0866	0.0752
1947Q1($T_1=176$)	1991Q1	0.0699°°	0.0818	0.0700	0.0697°	0.0706	0.0740	0.0748
1947Q1($T_1=188$)	1994Q1	0.0816	0.0943	0.0776°	0.0818	0.0814°°	0.0842	0.0887
1947Q1($T_1=200$)	1997Q1	0.0875°°	0.1006	0.0924	0.0870°	0.0885	0.0901	0.0923
1947Q1($T_1=212$)	2000Q1	0.0821°°	0.0951	0.0941	0.0774°	0.0849	0.0863	0.0859
1947Q1($T_1=224$)	2003Q1	0.0395	0.0520	0.0513	0.0397	0.0344°°	0.0339°	0.0512

Notes: (1) For comparison, all results are multiplied by 10. The estimation sample begins from 1947Q1, with T_1 observations. The prediction period begins from the quarter indicated in the second column.

(2) Seven methods are shown in Table 2: TVJMA, AICc in Cai & Tiwari (2000), SAICc, JMA, MMA, SAIC and SBIC estimators. The smaller the criteria, the better the method.

(3) The bandwidth used here is $2.34T_1^{-0.2}$, the same as that in the simulation study.

(4) ° and °° denote the best and the second best forecast among these seven methods, respectively.

References

- ALLEN, D. M. (1974). The relationship between variable selection and data agumentation and a method for prediction. *Technometrics* **16**, 125–127.
- ANDO, T. & LI, K.-C. (2014). A model-averaging approach for high-dimensional regression. *Journal of American Statistical Association* **109**, 254–265.
- ANDREWS, D. W. (1991). Asymptotic optimality of generalized C_L , cross-validation, and generalized cross-validation in regression with heteroskedastic errors. *Journal of Econometrics* **47**, 359–377.
- BAKER, M. & WURGLER, J. (2000). The equity share in new issues and aggregate stock returns. *Journal of Finance* **55**, 2219–2257.
- BLUNDELL, R., DUNCAN, A. & PENDAKUR, K. (1998). Semiparametric estimation and consumer demand. *Journal of Applied Econometrics* **13**, 435–461.
- BUCKLAND, S. T., BURNHAM, K. P. & AUGUSTIN, N. H. (1997). Model selection: an integral part of inference. *Biometrics* **53**, 603–618.
- CAI, Z. (2007). Trending time-varying coefficient time series models with serially correlated errors. *Journal of Econometrics* **136**, 163–188.
- CAI, Z. & TIWARI, R. C. (2000). Application of a local linear autoregressive model to bod time series. *Environmetrics* **11**, 341–350.
- CAMPBELL, J. Y. (1990). A variance decomposition for stock returns. *Economic Journal* **101**, 157–179.
- CAMPBELL, J. Y. & SHILLER, R. J. (1988). The dividend-price ratio and expectations of future dividends and discount factors. *Review of Financial Studies* **1**, 195–228.
- CAMPBELL, J. Y. & THOMPSON, S. B. (2008). Predicting excess stock returns out of sample: can anything beat the historical average? *Review of Financial Studies* **21**, 1509–1531.
- CHAN, K. S. & TONG, H. (1986). On estimating thresholds in autoregressive models. *Journal of Time Series Analysis* **7**, 179–190.
- CHEN, B. (2015). Modeling and testing smooth structural changes with endogenous regressors. *Journal of Econometrics* **185**, 196–215.
- CHEN, B. & HONG, Y. (2012). Testing for smooth structural changes in time series models via nonparametric regression. *Econometrica* **80**, 1157–1183.
- COCHRANE, J. H. (1996). A cross-sectional test of an investment-based asset pricing model. *Journal of Political Economy* **104**, 572–621.

- DAHLHAUS, R. (1996). On the Kullback-Leibler information divergence of locally stationary processes. *Stochastic Processes and Their Applications* **62**, 139–168.
- DAHLHAUS, R. (1997). Fitting time series models to nonstationary processes. *Annals of Statistics* **25**, 1–37.
- EPANECHNIKOV, V. A. (1969). Non-parametric estimation of a multivariate probability density. *Theory of Probability & Its Applications* **14**, 153–158.
- FAMA, E. F. & FRENCH, K. R. (1988). Dividend yields and expected stock returns. *Journal of Financial Economics* **22**, 3–25.
- GAO, Y., ZHANG, X. & WANG, S. (2016). Model averaging based on leave-subject-out cross-validation. *Journal of Econometrics* **192**, 139–151.
- GEISSER, S. (1975). The predictive sample reuse method with applications. *Journal of American Statistical Association* **70**, 320–328.
- GRANT, A. P. (2002). Time-varying estimates of the natural rate of unemployment: a revisitation of Okun’s law. *Quarterly Review of Economics and Finance* **42**, 95–113.
- HANSEN, B. E. (2001). The new econometrics of structural change: dating breaks in US labor productivity. *Journal of Economic Perspectives* **15**, 117–128.
- HANSEN, B. E. (2007). Least squares model averaging. *Econometrica* **75**, 1175–1189.
- HANSEN, B. E. (2008). Least-squares forecast averaging. *Journal of Econometrics* **146**, 342–350.
- HANSEN, B. E. & RACINE, J. S. (2012). Jackknife model averaging. *Journal of Econometrics* **167**, 38–46.
- HENDERSON, D. J. & PARMETER, C. F. (2016). Model averaging over nonparametric estimators. *Advances in Econometrics* **36**, 539–560.
- HJORT, N. L. & CLAESKENS, G. (2003). Frequentist model average estimators. *Journal of American Statistical Association* **98**, 879–899.
- HOETING, J. A., MADIGAN, D., RAFTERY, A. E. & VOLINSKY, C. T. (1999). Bayesian model averaging: a tutorial. *Statistical Science* **14**, 382–401.
- ING, C. K. & WEI, C. Z. (2003). On same-realization prediction in an infinite-order autoregressive process. *Journal of Multivariate Analysis* **85**, 130–155.
- JIN, S., SU, L. & ULLAH, A. (2014). Robustify financial time series forecasting with bagging. *Econometric Reviews* **33**, 575–605.
- LAMONT, O. (1998). Earnings and expected returns. *Journal of Finance* **53**, 1563–1587.

- LEHMANN, E. L. & CASELLA, G. (2006). *Theory of point estimation*. Springer Science & Business Media.
- LEWELLEN, J. (2004). Predicting returns with financial ratios. *Journal of Financial Economics* **74**, 209–235.
- LI, K. C. (1987). Asymptotic optimality for C_p , CL , cross-validation and generalized cross-validation: discrete index set. *Annals of Statistics* **15**, 958–975.
- LIN, C. F. J. & TERÄSVIRTA, T. (1994). Testing the constancy of regression parameters against continuous structural change. *Journal of Econometrics* **62**, 211–228.
- LIU, C.-A. (2015). Distribution theory of the least squares averaging estimator. *Journal of Econometrics* **186**, 142–159.
- LIU, Q. & OKUI, R. (2013). Heteroskedasticity-robust C_p model averaging. *Econometrics Journal* **16**, 462–473.
- LU, X. & SU, L. (2015). Jackknife model averaging for quantile regressions. *Journal of Econometrics* **188**, 40–58.
- PESARAN, M. H. & TIMMERMANN, A. (2007). Selection of estimation window in the presence of breaks. *Journal of Econometrics* **137**, 134–161.
- RAO, C. R. (1973). Linear statistical inference and its applications. *New York: Wiley* **2**.
- RAPACH, D. E., STRAUSS, J. K. & ZHOU, G. (2010). Out-of-sample equity premium prediction: Combination forecasts and links to the real economy. *Review of Financial Studies* **23**, 821–862.
- RAPACH, D. E. & ZHOU, G. (2013). Forecasting stock returns. *Handbook of Economic Forecasting* **2**, 328–383.
- ROBINSON, P. M. (1989). Nonparametric estimation of time-varying parameters. *Statistical Analysis and Forecasting of Economic Structural Change* **19**, 253–264.
- ROBINSON, P. M. (1991). Consistent nonparametric entropy-based testing. *Review of Economic Studies* **58**, 437–453.
- ROSSI, B. (2006). Are exchange rates really random walks? Some evidence robust to parameter instability. *Macroeconomic Dynamics* **10**, 20–38.
- ROSSI, B. & SEKHPOSYAN, T. (2011). Understanding models forecasting performance. *Journal of Econometrics* **164**, 158–172.
- ROTHMAN, P. (1998). Forecasting asymmetric unemployment rates. *Review of Economics and Statistics* **80**, 164–168.

- ROZEFF, M. S. (1984). Dividend yields are equity risk premiums. *Journal of Portfolio Management* **11**, 68–75.
- SHAO, J. (1997). An asymptotic theory for linear model selection. *Statistica Sinica* **7**, 221–242.
- STOCK, J. H. & WATSON, M. W. (1996). Evidence on structural instability in macroeconomic time series relations. *Journal of Business & Economic Statistics* **14**, 11–30.
- STOCK, J. H. & WATSON, M. W. (2002). Has the business cycle changed and why? *NBER Macroeconomics Annual* **17**, 159–218.
- STOCK, J. H. & WATSON, M. W. (2003). Forecasting output and inflation: the role of asset prices. *Journal of Economic Literature* **41**, 788–829.
- STOCK, J. H. & WATSON, M. W. (2005). Understanding changes in international business cycle dynamics. *Journal of European Economic Association* **3**, 968–1006.
- STOCK, J. H. & WATSON, M. W. (2007). Why has US inflation become harder to forecast? *Journal of Money, Credit and Banking* **39**, 3–33.
- STOCK, J. H. & WATSON, M. W. (2012). Generalized shrinkage methods for forecasting using many predictors. *Journal of Business & Economic Statistics* **30**, 481–493.
- STONE, M. (1974). Cross-validatory choice and assessment of statistical predictions. *Journal of Royal Statistical Society. Series B (Methodological)* **36**, 111–147.
- SUN, Y., HONG, Y. & WANG, S. (2018). Selection of an optimal rolling window in time-varying predictive regression. Manuscript, Cornell University .
- TWRDY, E. & BATISTA, M. (2016). Modeling of container throughput in northern adriatic ports over the period 1990–2013. *Journal of Transport Geography* **52**, 131–142.
- ULLAH, A., WAN, A. T., WANG, H., ZHANG, X. & ZOU, G. (2017). A semiparametric generalized ridge estimator and link with model averaging. *Econometric Reviews* **36**, 370–384.
- VOGT, M. (2012). Nonparametric regression for locally stationary time series. *Annals of Statistics* **40**, 2601–2633.
- WAN, A. T., ZHANG, X. & ZOU, G. (2010). Least squares model averaging by Mallows criterion. *Journal of Econometrics* **156**, 277–283.
- WANG, K. Q. (2003). Asset pricing with conditioning information: a new test. *Journal of Finance* **58**, 161–196.
- WELCH, I. & GOYAL, A. (2008). A comprehensive look at the empirical performance of equity premium prediction. *Review of Financial Studies* **21**, 1455–1508.

- YANG, Y. (2001). Adaptive regression by mixing. *Journal of American Statistical Association* **96**, 574–588.
- YUAN, Z. & YANG, Y. (2005). Combining linear regression models: when and how? *Journal of American Statistical Association* **100**, 1202–1214.
- ZHANG, X. (2015). Consistency of model averaging estimators. *Economics Letters* **130**, 120–123.
- ZHANG, X. & LIU, C.-A. (2018). Inference after model averaging in linear regression models. *Econometric Theory* , DOI:10.1017/S0266466618000269.
- ZHANG, X., WAN, A. T. & ZOU, G. (2013). Model averaging by jackknife criterion in models with dependent data. *Journal of Econometrics* **174**, 82–94.
- ZHU, R., ZHANG, X., WAN, A. T. & ZOU, G. (2017). Kernel averaging estimators. *Working Paper* .

Appendix

1 Appendix A.1

Proof of Theorem 1. First we show

$$\tilde{L}_{t,T}(\hat{\mathbf{w}}_t) / \inf_{\mathbf{w} \in \mathcal{H}_T} \tilde{L}_{t,T}(\mathbf{w}) \xrightarrow{p} 1 \quad (\text{A.1})$$

with Assumptions 1 - 7.

Note that

$$\begin{aligned} \text{CV}_{t,T}(\mathbf{w}) &= \tilde{L}_{t,T}(\mathbf{w}) + \boldsymbol{\varepsilon}' \mathbf{K}_t \boldsymbol{\varepsilon} + 2\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu} - 2\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} \\ &= \tilde{L}_{t,T}(\mathbf{w}) + \boldsymbol{\varepsilon}' \mathbf{K}_t \boldsymbol{\varepsilon} + 2\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu} - 2(\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega})) - 2\text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega}). \end{aligned}$$

With Assumption 7, (A.1) is valid if the following results hold: as $T \rightarrow \infty$,

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0, \quad (\text{A.2})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega})}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0, \quad (\text{A.3})$$

and

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\tilde{L}_{t,T}(\mathbf{w})}{\tilde{R}_{t,T}(\mathbf{w})} - 1 \right| \xrightarrow{p} 0. \quad (\text{A.4})$$

(A.2)-(A.4) will be verified later.

We next show that

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{L_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right| \xrightarrow{p} 0. \quad (\text{A.5})$$

It is straightforward to obtain that

$$\begin{aligned} &L_{t,T}(\mathbf{w}) - R_{t,T}(\mathbf{w}) \\ &= \boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu} + \boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} - 2\boldsymbol{\mu}' \mathbf{A}(\mathbf{w})' \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} \\ &\quad - [\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu} + \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega})] \\ &= \boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}(\mathbf{w})' \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega}) - 2\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} \\ &= D_1(\mathbf{w}) - D_2(\mathbf{w}), \end{aligned}$$

where $D_1(\mathbf{w}) = \boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}(\mathbf{w})' \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega})$ and $D_2(\mathbf{w}) = 2\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}$.

Thus, to prove (A.5), we only need to verify the following two results:

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{D_1(\mathbf{w})}{R_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0. \quad (\text{A.6})$$

and

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{D_2(\mathbf{w})}{R_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0. \quad (\text{A.7})$$

Because of the normality assumption, $\mathbf{\Omega}^{-1/2}\boldsymbol{\varepsilon}$ is a vector of independent variables, and this allows us to use Theorem 2 in Whittle (1960). Thus, based on the Chebyshev's inequality, Theorem 2 in Whittle (1960) and Assumption 1, we have, when \mathbf{X} is nonstochastic, for any $\delta > 0$,

$$\begin{aligned} & \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \mathbf{\Omega})}{R_{t,T}(\mathbf{w})} \right| > \delta \right) \\ & \leq \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \mathbf{\Omega})| > \delta \xi_{t,T} \right) \\ & \leq \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m |\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \mathbf{\Omega})| > \delta \xi_{t,T} \right) \\ & \leq \Pr \left(\max_{1 \leq k \leq M_T} \max_{1 \leq m \leq M_T} |\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \mathbf{\Omega})| > \delta \xi_{t,T} \right) \\ & = \Pr \{ (|\boldsymbol{\varepsilon}' \mathbf{P}'_1 \mathbf{K}_t \mathbf{P}_1 \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_1 \mathbf{K}_t \mathbf{P}_1 \mathbf{\Omega})| > \delta \xi_{t,T}) \cup \\ & \quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_1 \mathbf{K}_t \mathbf{P}_2 \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_1 \mathbf{K}_t \mathbf{P}_2 \mathbf{\Omega})| > \delta \xi_{t,T}) \cup \dots \cup \\ & \quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_1 \mathbf{K}_t \mathbf{P}_{M_T} \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_1 \mathbf{K}_t \mathbf{P}_{M_T} \mathbf{\Omega})| > \delta \xi_{t,T}) \cup \\ & \quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_2 \mathbf{K}_t \mathbf{P}_1 \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_2 \mathbf{K}_t \mathbf{P}_1 \mathbf{\Omega})| > \delta \xi_{t,T}) \cup \dots \cup \\ & \quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_{M_T} \mathbf{K}_t \mathbf{P}_{M_T} \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_{M_T} \mathbf{K}_t \mathbf{P}_{M_T} \mathbf{\Omega})| > \delta \xi_{t,T}) \} \\ & \leq \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \Pr (|\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \mathbf{\Omega})| > \delta \xi_{t,T}) \\ & \leq \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \mathbb{E} \left[\frac{(\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \mathbf{\Omega}))^{2G}}{\delta^{2G} \xi_{t,T}^{2G}} \right] \\ & \leq C_4 \delta^{-2G} \xi_{t,T}^{-2G} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \left[\text{tr} \left((\mathbf{P}'_k \mathbf{K}_t \mathbf{P}_m \mathbf{\Omega})^2 \right) \right]^G \\ & \leq C_4 \delta^{-2G} \xi_{t,T}^{-2G} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} [\zeta(\mathbf{P}_m \mathbf{\Omega} \mathbf{P}'_m) k_{\max} \text{tr}(\mathbf{P}'_k \mathbf{K}_t \mathbf{P}_k \mathbf{\Omega})]^G \\ & \leq C'_4 \delta^{-2G} \xi_{t,T}^{-2G} M_T \sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0))^G, \end{aligned} \quad (\text{A.8})$$

where \mathbf{w}_m^0 defines the weight with the m -th element 1 and others 0, and C_4 and C'_4 are some constants. Thus, given Assumption 6, (A.6) with non-stochastic \mathbf{X} is verified.

To prove (A.7), we have $\|\boldsymbol{\mu}' \mathbf{A}'_k \mathbf{K}_t \mathbf{P}_m\|^2 \leq \zeta^2(\mathbf{P}_m) k_{\max} \boldsymbol{\mu}' \mathbf{A}'_k \mathbf{K}_t \mathbf{A}_k \boldsymbol{\mu} \leq C \boldsymbol{\mu}' \mathbf{A}'_k \mathbf{K}_t \mathbf{A}_k \boldsymbol{\mu}$, $1 \leq k \leq M_T$, $1 \leq m \leq M_T$ with Assumption 3. Based on the Chebyshev's inequality,

Theorem 2 in Whittle (1960) and Assumption 1, for any $\delta > 0$, we have

$$\begin{aligned}
& \Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}_t(\mathbf{w}) \boldsymbol{\varepsilon}}{R_{t,T}(\mathbf{w})} \right| > \delta \right\} \\
& \leq \Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \sum_{1 \leq m \leq M_T} \sum_{1 \leq k \leq M_T} w^k w^m |\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_k)' \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon}| > \delta \xi_{t,T} \right\} \\
& \leq \Pr \left\{ \max_{1 \leq k \leq M_T} \max_{1 \leq m \leq M_T} |\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_k)' \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon}| > \delta \xi_{t,T} \right\} \\
& = \Pr \{ (|\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_1)' \mathbf{K}_t \mathbf{P}_1 \boldsymbol{\varepsilon}| > \delta \xi_{t,T}) \cup \\
& \quad (|\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_1)' \mathbf{K}_t \mathbf{P}_2 \boldsymbol{\varepsilon}| > \delta \xi_{t,T}) \cup \dots \cup \\
& \quad (|\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_1)' \mathbf{K}_t \mathbf{P}_{M_T} \boldsymbol{\varepsilon}| > \delta \xi_{t,T}) \cup \\
& \quad (|\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_2)' \mathbf{K}_t \mathbf{P}_1 \boldsymbol{\varepsilon}| > \delta \xi_{t,T}) \cup \dots \cup \\
& \quad (|\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_{M_T})' \mathbf{K}_t \mathbf{P}_{M_T} \boldsymbol{\varepsilon}| > \delta \xi_{t,T}) \} \\
& \leq \sum_{1 \leq k \leq M_T} \sum_{1 \leq m \leq M_T} \mathbb{E} \left[\frac{(\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_k)' \mathbf{K}_t \mathbf{P}_m \boldsymbol{\varepsilon})^{2G}}{\delta^{2G} \xi_{t,T}^{2G}} \right] \\
& \leq C_5' \delta^{-2G} \xi_{t,T}^{-2G} \sum_{1 \leq k \leq M_T} \sum_{1 \leq m \leq M_T} \|\boldsymbol{\mu}'(\mathbf{I}_T - \mathbf{P}_k)' \mathbf{K}_t \mathbf{P}_m\|^{2G} \\
& \leq C_5' \delta^{-2G} \xi_{t,T}^{-2G} \sum_{1 \leq k \leq M_T} \sum_{1 \leq m \leq M_T} (\boldsymbol{\mu}' \mathbf{A}'_k \mathbf{K}_t \mathbf{A}_m \boldsymbol{\mu})^G \\
& \leq C_5' \delta^{-2G} \xi_{t,T}^{-2G} M_T \sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0))^G, \tag{A.9}
\end{aligned}$$

where C_5 and C_5' are some constants. Given Assumption 6, (A.7) with non-stochastic \mathbf{X} is verified.

Besides, for the case of random \mathbf{X} , based on the dominated convergence theorem, Assumption 5 and (A.2), the results in (A.6) and (A.7) are obtained from (A.8) and (A.9), respectively.

Define

$$V_{t,T}(\widehat{\mathbf{w}}_t) = \boldsymbol{\mu}' \mathbf{A}'(\widehat{\mathbf{w}}_t) \mathbf{K}_t \mathbf{A}(\widehat{\mathbf{w}}_t) \boldsymbol{\mu} + \text{tr}(\mathbf{P}'(\widehat{\mathbf{w}}_t) \mathbf{K}_t \mathbf{P}(\widehat{\mathbf{w}}_t) \boldsymbol{\Omega}),$$

and

$$\widetilde{V}_n(\widehat{\mathbf{w}}_t) = \boldsymbol{\mu}' \widetilde{\mathbf{A}}'(\widehat{\mathbf{w}}_t) \mathbf{K}_t \widetilde{\mathbf{A}}(\widehat{\mathbf{w}}_t) \boldsymbol{\mu} + \text{tr}(\widetilde{\mathbf{P}}'(\widehat{\mathbf{w}}_t) \mathbf{K}_t \widetilde{\mathbf{P}}(\widehat{\mathbf{w}}_t) \boldsymbol{\Omega}).$$

Next, we follow the spirit in Zhang et al. (2013) to obtain

$$\begin{aligned}
& \frac{L_{t,T}(\widehat{\mathbf{w}}_t)}{\inf_{\mathbf{w} \in \mathcal{H}_T} L_{t,T}(\mathbf{w})} - 1 = \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{L_{t,T}(\widehat{\mathbf{w}}_t)}{L_{t,T}(\mathbf{w})} - 1 \right) \\
& = \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{L_{t,T}(\widehat{\mathbf{w}}_t)}{V_{t,T}(\widehat{\mathbf{w}}_t)} \frac{R_{t,T}(\mathbf{w})}{L_{t,T}(\mathbf{w})} \frac{\widetilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} \frac{V_{t,T}(\widehat{\mathbf{w}}_t)}{\widetilde{V}_{t,T}(\widehat{\mathbf{w}}_t)} \right)
\end{aligned}$$

$$\begin{aligned}
& \times \frac{\tilde{V}_{t,T}(\hat{\mathbf{w}}_t)}{\tilde{L}_{t,T}(\hat{\mathbf{w}}_t)} \frac{\tilde{L}_{t,T}(\mathbf{w})}{\tilde{R}_{t,T}(\mathbf{w})} \frac{\tilde{L}_{t,T}(\hat{\mathbf{w}}_t)}{\tilde{L}_{t,T}(\mathbf{w})} - 1 \Bigg) \\
& \leq \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{L_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} \right) \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{R_{t,T}(\mathbf{w})}{L_{t,T}(\mathbf{w})} \right) \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} \right) \\
& \quad \times \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{R_{t,T}(\mathbf{w})}{\tilde{R}_{t,T}(\mathbf{w})} \right) \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{\tilde{R}_{t,T}(\mathbf{w})}{\tilde{L}_{t,T}(\mathbf{w})} \right) \sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{\tilde{L}_{t,T}(\mathbf{w})}{\tilde{R}_{t,T}(\mathbf{w})} \right) \frac{\tilde{L}_{t,T}(\hat{\mathbf{w}}_t)}{\inf_{\mathbf{w} \in \mathcal{H}_T} \tilde{L}_{t,T}(\mathbf{w})} - 1.
\end{aligned}$$

Thus, given (A.2)-(A.4), Theorem 1 holds.

Next, we only need to verify (A.2)-(A.4). Denote $\tilde{\xi}_{t,T} = \inf_{\mathbf{w} \in \mathcal{H}_T} \tilde{R}_{t,T}(\mathbf{w})$. If $\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right| \leq 1$, we have

$$\begin{aligned}
& \tilde{\xi}_{t,T}^{2G} \left[\sum_{m=1}^{M_T} (\tilde{R}_{t,T}(\mathbf{w}_m^0))^G \right]^{-1} \\
& = \left[\inf_{\mathbf{w} \in \mathcal{H}_{t,T}} \left(R_{t,T}(\mathbf{w}) \frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} \right) \right]^{2G} \left[\sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0))^G \left(\frac{\tilde{R}_{t,T}(\mathbf{w}_m^0)}{R_{t,T}(\mathbf{w}_m^0)} \right)^G \right]^{-1} \\
& \geq \xi_{t,T}^{2G} \left[\inf_{\mathbf{w} \in \mathcal{H}_{t,T}} \frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} \right]^{2G} \left[\max_{1 \leq m \leq M_T} \frac{\tilde{R}_{t,T}(\mathbf{w}_m^0)}{R_{t,T}(\mathbf{w}_m^0)} \right]^{-G} \left[\sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0))^G \right]^{-1} \\
& \geq \xi_{t,T}^{2G} \left[1 + \inf_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right) \right]^{2G} \left[\sup_{\mathbf{w} \in \mathcal{H}_T} \left(\frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right) + 1 \right]^{-G} \left[\sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0))^G \right]^{-1} \\
& \geq \xi_{t,T}^{2G} \left[1 - \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right| \right]^{2G} \left[\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right| + 1 \right]^{-G} \\
& \quad \times \left[\sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0))^G \right]^{-1}. \tag{A.10}
\end{aligned}$$

Given Assumptions 5 and 6, we obtain the following result from (A.10):

$$M_T \tilde{\xi}_{t,T}^{-2G} \sum_{m=1}^{M_T} (\tilde{R}_{t,T}(\mathbf{w}_m^0))^G \xrightarrow{a.s.} 0. \tag{A.11}$$

To prove (A.2), by using the Chebyshev inequality, Theorem 2 of Whittle (1960) and Assumptions 1-2, we have, for any $\delta > 0$:

$$\begin{aligned}
& \Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}}{\tilde{R}_{t,T}(\mathbf{w})} \right| > \delta \right\} \\
& \leq \Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \sum_{m=1}^M w^m |\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T} \right\} \\
& = \Pr \left\{ \max_{1 \leq m \leq M_T} |\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T} \right\}
\end{aligned}$$

$$\begin{aligned}
&= \Pr \left\{ (|\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(w_1^0) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T}) \cup \dots \cup (|\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(w_{M_T}^0) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T}) \right\} \\
&\leq \sum_{m=1}^{M_T} \Pr \left\{ |\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T} \right\} \\
&\leq \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \sum_{m=1}^{M_T} E(\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\varepsilon})^{2G} \\
&\leq C_1 \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \sum_{m=1}^{M_T} [\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\Omega} \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}]^G \\
&\leq C_1 \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \delta^G(\boldsymbol{\Omega}) \sum_{m=1}^{M_T} (\text{tr}(\mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu} \boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}_m^0) \mathbf{K}_t)) \\
&\leq C_1 \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \delta^G(\boldsymbol{\Omega}) k_{\max} \sum_{m=1}^{M_T} [\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}_m^0) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}]^G,
\end{aligned}$$

where C_1 is a positive constant. Then from (A.11) and Assumptions 2 and 4, we have $\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0$.

Next, we verify (A.3).

$$\begin{aligned}
&\Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega})}{\tilde{R}_{t,T}(\mathbf{w})} \right| > \delta \right\} \\
&\leq \sum_{m=1}^{M_T} \Pr \left\{ \sup \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\Omega})}{\tilde{R}_{t,T}(\mathbf{w}_m^0)} \right| > \delta \right\} \\
&\leq \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \sum_{m=1}^{M_T} \mathbb{E}[|\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\Omega})|]^{2G} \\
&\leq C_2 \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \sum_{m=1}^{M_T} \left[\text{tr} \left(\boldsymbol{\Omega}^{\frac{1}{2}} \tilde{\mathbf{P}}(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\Omega} \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\Omega}^{\frac{1}{2}} \right) \right] \\
&\leq C_2 \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \zeta(\boldsymbol{\Omega}) \sum_{m=1}^{M_T} \left[\text{tr} \left(\tilde{\mathbf{P}}(\mathbf{w}_m^0) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\Omega} \right) \right] \\
&\leq C_2 \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \zeta(\boldsymbol{\Omega}) k_{\max} \sum_{m=1}^{M_T} \left[\text{tr} \left(\tilde{\mathbf{P}}'(\mathbf{w}_m^0) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}_m^0) \boldsymbol{\Omega} \right) \right],
\end{aligned}$$

where C_2 is a positive constant. Then from Assumptions 2, 4 and (A.11), we obtain that $\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega})}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0$.

Finally, we will complete the proof of Theorem 1 by verifying (A.4). Note that

$$\begin{aligned}
&|\tilde{L}_{t,T}(\mathbf{w}) - \tilde{R}_{t,T}(\mathbf{w})| \\
&= |\mathbf{e}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \mathbf{e} - \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega}) - 2\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}|.
\end{aligned}$$

We only need to prove the following results:

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0 \quad (\text{A.12})$$

and

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega})}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0. \quad (\text{A.13})$$

For (A.12), we have that

$$\begin{aligned} & \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}}{\tilde{R}_{t,T}(\mathbf{w})} \right| > \delta \right) \\ & \leq \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T} \right) \\ & \leq \sum_{m=1}^{M_T} \Pr \left(|\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}_m^0) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}| > \delta \tilde{\xi}_{t,T} \right) \\ & \leq \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \sum_{m=1}^{M_T} \mathbb{E} |\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}_m^0) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu}|^{2G} \\ & \leq \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} \max_{1 \leq m \leq M_T} \tilde{\mathbf{P}}'(\mathbf{w}_m^0) \left(\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\Omega} \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu} \right)^G \\ & \leq \delta^{-2G} \tilde{\xi}_{t,T}^{-2G} C k_{\max} \zeta(\boldsymbol{\Omega}) \left(\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}_m^0) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}_m^0) \boldsymbol{\mu} \right)^G. \end{aligned}$$

Then from (A.11) and Assumptions 2 and 4, $\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0$. The proof of (A.13) is similar to that of (A.12), and thus the proof is omitted. It follows that the proof of (A.3) is completed. \blacksquare

2 Appendix A.2

Proof of Theorem 2. Denote the maximum singular values of matrixes $B_i, i = 1, 2$, by $\zeta(B_1)$ and $\zeta(B_2)$. It is known that for any square matrices B_1 and B_2 with identical dimensions, the following simple inequalities are obtained:

$$\zeta(B_1 B_2) \leq \zeta(B_1) \zeta(B_2), \quad (\text{A.14})$$

and

$$\zeta(B_1 + B_2) \leq \zeta(B_1) + \zeta(B_2). \quad (\text{A.15})$$

See more discussions of these inequalities in proof of Theorem 5.2 in Li (1987) and proof of Theorem 2.2 in Zhang et al. (2013).

Let $\tilde{h} = \frac{h^*}{1-h^*}$. By Assumption 12, we have $h^* = O(T^{-1}h^{-1})$ a.s., $\tilde{h} = O(T^{-1}h^{-1})$ a.s., and then given Assumption 9, we have

$$h^* \rightarrow 0 \text{ and } \tilde{h} \rightarrow 0, \text{ a.s..} \quad (\text{A.16})$$

Let \mathbf{Q}_m be a $T \times T$ diagonal matrix with the (t, t) -th element $\frac{h_{tt}^m}{1-h_{tt}^m}$. Then it is easy to obtain that $\tilde{\mathbf{P}}_m = \mathbf{P}_m - \mathbf{Q}_m \mathbf{A}_m$, and $\tilde{\mathbf{A}}_m = \mathbf{A}_m + \mathbf{Q}_m \mathbf{A}_m$. Denote $\mathbf{Q}(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \mathbf{Q}_m$, $\mathbf{T}_m = \mathbf{Q}_m \mathbf{P}_m$ and $\mathbf{T}(\mathbf{w}) = \sum_{m=1}^{M_T} w^m \mathbf{T}_m$. To prove Theorem 2, we only need to verify the following results:

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{L_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right| \xrightarrow{p} 0, \quad (\text{A.17})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\tilde{R}_{t,T}(\mathbf{w})}{R_{t,T}(\mathbf{w})} - 1 \right| \xrightarrow{p} 0, \quad (\text{A.18})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\tilde{L}_{t,T}(\mathbf{w})}{\tilde{R}_{t,T}(\mathbf{w})} - 1 \right| \xrightarrow{p} 0, \quad (\text{A.19})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\mu}' \tilde{\mathbf{A}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0, \quad (\text{A.20})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \xrightarrow{p} 0. \quad (\text{A.21})$$

Intuitively, (A.17), (A.19) and (A.20) are similar to (A.5), (A.4) and (A.2) in proof of Theorem 1. To prove (A.17), we have:

$$\begin{aligned} & \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|L_{t,T}(\mathbf{w}) - R_{t,T}(\mathbf{w})|}{R_{t,T}(\mathbf{w})} \\ & \leq \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}|}{R_{t,T}(\mathbf{w})} + \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}|}{R_{t,T}(\mathbf{w})} \\ & + \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{\text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega})}{R_{t,T}(\mathbf{w})} \\ & \equiv \Delta_{11} + \Delta_{12} + \Delta_{13}. \end{aligned}$$

We next show that $\Delta_{11} = o_p(1)$, $\Delta_{12} = o_p(1)$ and $\Delta_{13} = o_p(1)$. For Δ_{11} , we have

$$|\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}| \leq k_{\max} \boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}. \quad (\text{A.22})$$

For Δ_{12} , we have

$$\frac{|\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}|}{R_{t,T}(\mathbf{w})}$$

$$\begin{aligned}
&\leq \left[k_{\max} \frac{\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}}{R_{t,T}(\mathbf{w})} \frac{\boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu}}{R_{t,T}(\mathbf{w})} \right]^{\frac{1}{2}} \\
&\leq \left[k_{\max} \frac{\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}}{R_{t,T}(\mathbf{w})} \right]^{\frac{1}{2}}.
\end{aligned} \tag{A.23}$$

Moreover, it is straightforward to obtain that

$$\begin{aligned}
&\Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon}}{R_{t,T}(\mathbf{w})} - \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega}) \right| > \delta \right) \\
&\leq \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega})| > \delta \xi_{t,T} \right) \\
&\leq \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m |\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{P}_m \boldsymbol{\Omega})| > \delta \xi_{t,T} \right) \\
&\leq \Pr \left(\max_{1 \leq k \leq M_T} \max_{1 \leq m \leq M_T} |\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{P}_m \boldsymbol{\Omega})| > \delta \xi_{t,T} \right) \\
&= \Pr \{ (|\boldsymbol{\varepsilon}' \mathbf{P}'_1 \mathbf{P}_1 \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_1 \mathbf{P}_1 \boldsymbol{\Omega})| > \delta \xi_{t,T}) \cup \\
&\quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_1 \mathbf{P}_2 \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_1 \mathbf{P}_2 \boldsymbol{\Omega})| > \delta \xi_{t,T}) \cup \dots \cup \\
&\quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_1 \mathbf{P}_{M_T} \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_1 \mathbf{P}_{M_T} \boldsymbol{\Omega})| > \delta \xi_{t,T}) \cup \\
&\quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_2 \mathbf{P}_1 \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_2 \mathbf{P}_1 \boldsymbol{\Omega})| > \delta \xi_{t,T}) \cup \dots \cup \\
&\quad (|\boldsymbol{\varepsilon}' \mathbf{P}'_{M_T} \mathbf{P}_{M_T} \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_{M_T} \mathbf{P}_{M_T} \boldsymbol{\Omega})| > \delta \xi_{t,T}) \} \\
&\leq \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \Pr (|\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{P}_m \boldsymbol{\Omega})| > \delta \xi_{t,T}) \\
&\leq \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \mathbb{E} \left[\frac{(\boldsymbol{\varepsilon}' \mathbf{P}'_k \mathbf{P}_m \boldsymbol{\varepsilon} - \text{tr}(\mathbf{P}'_k \mathbf{P}_m \boldsymbol{\Omega}))^2}{\delta^2 \xi_{t,T}^{2G}} \right] \\
&\leq C \delta^{-2} \xi_{t,T}^{-2} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \text{tr} \left((\mathbf{P}'_k \mathbf{P}_m \boldsymbol{\Omega})^2 \right) \\
&\leq C' \delta^{-2} \xi_{t,T}^{-2} M_T \sum_{m=1}^{M_T} (R_{t,T}(\mathbf{w}_m^0)),
\end{aligned} \tag{A.24}$$

where C and C' are some constants. Then given Assumptions 2 and 10, we have

$$\begin{aligned}
&\Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \frac{\text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega})}{R_{t,T}(\mathbf{w})} \geq \delta \right) \\
&\leq \Pr \left(\sup_{\mathbf{w} \in \mathcal{H}_T} \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \boldsymbol{\Omega}) \geq \xi_{t,T} \delta \right) \\
&\leq \Pr \left(\zeta(\boldsymbol{\Omega}) \sup_{\mathbf{w} \in \mathcal{H}_T} \text{tr}(\mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w})) \xi_{t,T}^{-1} \geq \delta \right) \\
&\rightarrow 0.
\end{aligned}$$

Combining this result with (A.22), (A.23) and (A.24), we have $\Delta_{11} = o_p(1)$ and $\Delta_{12} = o_p(1)$.

For Δ_{13} , we have:

$$\begin{aligned}
& \text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{K}_t\mathbf{P}(\mathbf{w})\boldsymbol{\Omega}) \\
& \leq k_{\max}\text{tr}[\mathbf{P}(\mathbf{w})\boldsymbol{\Omega}\mathbf{P}'(\mathbf{w})] \\
& \leq k_{\max}\zeta(\boldsymbol{\Omega})\text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{P}(\mathbf{w})).
\end{aligned}$$

Then given Assumptions 2 and 10, we have

$$\begin{aligned}
& \Pr\left(\sup_{\mathbf{w} \in \mathcal{H}_T} \frac{\text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{K}_t\mathbf{P}(\mathbf{w})\boldsymbol{\Omega})}{R_{t,T}(\mathbf{w})} \geq \delta\right) \\
& \leq \Pr\left(\sup_{\mathbf{w} \in \mathcal{H}_T} \text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{K}_t\mathbf{P}(\mathbf{w})\boldsymbol{\Omega}) \geq \xi_{t,T}\delta\right) \\
& \leq \Pr\left(k_{\max}\zeta(\boldsymbol{\Omega}) \sup_{\mathbf{w} \in \mathcal{H}_T} \text{tr}(\mathbf{P}'(\mathbf{w})\mathbf{P}(\mathbf{w}))\xi_{t,T}^{-1} \geq \delta\right) \\
& \rightarrow 0,
\end{aligned}$$

then $\Delta_{13} = o_p(1)$. Thus, the proof of (A.17) is completed.

To prove (A.18), we have:

$$\begin{aligned}
\sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\tilde{R}_{t,T}(\mathbf{w}) - R_{t,T}(\mathbf{w})|}{R_{t,T}(\mathbf{w})} & \leq \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\boldsymbol{\mu}'\tilde{\mathbf{A}}'(\mathbf{w})\mathbf{K}_t\tilde{\mathbf{A}}(\mathbf{w})\boldsymbol{\mu} - \boldsymbol{\mu}'\mathbf{A}'(\mathbf{w})\mathbf{K}_t\mathbf{A}(\mathbf{w})\boldsymbol{\mu}|}{R_{t,T}(\mathbf{w})} \\
& \quad + \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\text{tr}(\tilde{\mathbf{P}}(\mathbf{w})\mathbf{K}_t\tilde{\mathbf{P}}'(\mathbf{w})\boldsymbol{\Omega})|}{R_{t,T}(\mathbf{w})} + \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\text{tr}(\mathbf{P}(\mathbf{w})\mathbf{K}_t\mathbf{P}'(\mathbf{w})\boldsymbol{\Omega})|}{R_{t,T}(\mathbf{w})} \\
& \equiv \Delta_{21} + \Delta_{22} + \Delta_{23}.
\end{aligned}$$

Because the matrix $\mathbf{P}'(\mathbf{w})\boldsymbol{\Omega}\mathbf{P}(\mathbf{w})$ is symmetric, $\Delta_{23} = \Delta_{13} = o_p(1)$. Then we only need to verify $\Delta_{21} = o_p(1)$ and $\Delta_{22} = o_p(1)$. For Δ_{22} ,

$$\begin{aligned}
& \text{tr}\left(\tilde{\mathbf{P}}(\mathbf{w})\mathbf{K}_t\tilde{\mathbf{P}}'(\mathbf{w})\boldsymbol{\Omega}\right) \\
& \leq \text{tr}(\mathbf{P}(\mathbf{w})\mathbf{K}_t\mathbf{P}'(\mathbf{w})\boldsymbol{\Omega}) + \text{tr}(\mathbf{Q}(\mathbf{w})\mathbf{K}_t\mathbf{Q}'(\mathbf{w})\boldsymbol{\Omega}) \\
& \quad + \text{tr}(\mathbf{T}(\mathbf{w})\mathbf{K}_t\mathbf{T}'(\mathbf{w})\boldsymbol{\Omega}) + 2|\text{tr}(\mathbf{Q}(\mathbf{w})\mathbf{K}_t\mathbf{P}'(\mathbf{w})\boldsymbol{\Omega})| \\
& \quad + 2|\text{tr}(\mathbf{Q}(\mathbf{w})\mathbf{K}_t\mathbf{T}'(\mathbf{w})\boldsymbol{\Omega})| + 2|\text{tr}(\mathbf{T}(\mathbf{w})\mathbf{K}_t\mathbf{P}'(\mathbf{w})\boldsymbol{\Omega})|. \tag{A.25}
\end{aligned}$$

Given $\text{rank}(\mathbf{P}_m) \leq T$, $h^* = O(T^{-1}h^{-1})$ a.s. and $\tilde{h} = O(T^{-1}h^{-1})$ a.s., terms on the right side of (A.25) follow that

$$\begin{aligned}
|\text{tr}(\mathbf{Q}(\mathbf{w})\mathbf{K}_t\mathbf{Q}'(\mathbf{w})\boldsymbol{\Omega})| & \leq \zeta(\mathbf{Q}(\mathbf{w}))k_{\max}|\text{tr}(\mathbf{Q}(\mathbf{w})\boldsymbol{\Omega})| \\
& \leq \tilde{h}k_{\max}|\text{tr}(\mathbf{Q}(\mathbf{w})\boldsymbol{\Omega})| \\
& \leq \frac{\tilde{h}}{1-h^*}\zeta(\boldsymbol{\Omega})k_{\max}\sum_{m=1}^{M_T}w^m|\text{tr}(\mathbf{P}_m)|
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{\tilde{h}}{1-h^*} \zeta(\boldsymbol{\Omega}) k_{\max} \max_{1 \leq m \leq M_T} \text{rank}(\mathbf{P}_m) \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) \\
&= O\left(h^{-1} \zeta(\boldsymbol{\Omega}) k_{\max} \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m)\right)
\end{aligned}$$

and

$$\begin{aligned}
|\text{tr}(\mathbf{Q}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \boldsymbol{\Omega})| &\leq \text{rank}(\mathbf{P}(\mathbf{w})) \zeta(\mathbf{Q}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \boldsymbol{\Omega}) \\
&\leq \text{rank}(\mathbf{P}(\mathbf{w})) \tilde{h} k_{\max} \zeta(\boldsymbol{\Omega}) \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) \\
&= O\left(h^{-1} \zeta(\boldsymbol{\Omega}) k_{\max} \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m)\right).
\end{aligned}$$

Similarly, for $\mathbf{T}(\mathbf{w})$ we have

$$\begin{aligned}
|\text{tr}(\mathbf{T}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \boldsymbol{\Omega})| &\leq \text{rank}(\mathbf{P}(\mathbf{w})) \zeta(\mathbf{T}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \boldsymbol{\Omega}) \\
&\leq \text{rank}(\mathbf{P}(\mathbf{w})) k_{\max} \zeta(\mathbf{T}(\mathbf{w})) \zeta(\mathbf{P}'(\mathbf{w})) \zeta(\boldsymbol{\Omega}) \\
&\leq \text{rank}(\mathbf{P}(\mathbf{w})) k_{\max} \zeta\left(\sum_{m=1}^{M_T} w^m \mathbf{P}_m \mathbf{Q}_m\right) \zeta(\mathbf{P}'(\mathbf{w})) \zeta(\boldsymbol{\Omega}) \\
&\leq \tilde{h} k_{\max} \zeta(\boldsymbol{\Omega}) \text{rank}(\mathbf{P}(\mathbf{w})) \max_{1 \leq m \leq M_T} \zeta^2(\mathbf{P}_m), \\
&= O\left(h^{-1} \zeta(\boldsymbol{\Omega}) k_{\max} \max_{1 \leq m \leq M_T} \zeta^2(\mathbf{P}_m)\right),
\end{aligned}$$

$$\begin{aligned}
|\text{tr}(\mathbf{T}(\mathbf{w}) \mathbf{K}_t \mathbf{Q}'(\mathbf{w}) \boldsymbol{\Omega})| &\leq \zeta(\boldsymbol{\Omega}) \tilde{h}^2 k_{\max} \text{rank}(\mathbf{P}(\mathbf{w})) \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m), \\
&= O\left(T^{-1} h^{-2} \zeta(\boldsymbol{\Omega}) k_{\max} \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m)\right),
\end{aligned}$$

and

$$\begin{aligned}
|\text{tr}(\mathbf{T}(\mathbf{w}) \mathbf{K}_t \mathbf{T}'(\mathbf{w}) \boldsymbol{\Omega})| &\leq \zeta(\boldsymbol{\Omega}) \tilde{h}^2 k_{\max} \text{rank}(\mathbf{P}(\mathbf{w})) \max_{1 \leq m, k \leq M_T} \zeta^2(\mathbf{P}_m) \\
&= O\left(T^{-1} h^{-2} \zeta(\boldsymbol{\Omega}) k_{\max} \max_{1 \leq m \leq M_T} \zeta^2(\mathbf{P}_m)\right).
\end{aligned}$$

Then given Assumption 12, we have $\Pr\left\{\sup_{\mathbf{w} \in \mathcal{H}_T} \left|\frac{\text{tr}(\tilde{\mathbf{P}}(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega})}{R_{t,T}(\mathbf{w})}\right| > \delta\right\} \rightarrow 0$. Thus, $\Delta_{22} = o_p(1)$.

For Δ_{21} , we have

$$\begin{aligned}
&\left|\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu} - \boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu}\right| / R_{t,T}(\mathbf{w}) \\
&= \left|\sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \left\{w^k w^m \left(\boldsymbol{\mu}' \tilde{\mathbf{A}}'_k \mathbf{K}_t \tilde{\mathbf{A}}_m \boldsymbol{\mu} - \boldsymbol{\mu}' \mathbf{A}'_k \mathbf{K}_t \mathbf{A}_m \boldsymbol{\mu}\right)\right\}\right| / R_{t,T}(\mathbf{w})
\end{aligned}$$

$$\begin{aligned}
&= \left| \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} \left\{ w^k w^m [\boldsymbol{\mu}'(\mathbf{A}_k + \mathbf{Q}_k \mathbf{A}_k)' \mathbf{K}_t (\mathbf{A}_m + \mathbf{Q}_m \mathbf{A}_m) \boldsymbol{\mu} - \boldsymbol{\mu}' \mathbf{A}_k' \mathbf{K}_t \mathbf{A}_m \boldsymbol{\mu}] \right\} \right| / R_{t,T}(\mathbf{w}) \\
&= \left| \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m [\boldsymbol{\mu}' \mathbf{A}_k' \mathbf{Q}_k' \mathbf{K}_t \mathbf{Q}_m \boldsymbol{\mu}] + 2 \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\mu}' \mathbf{A}_k' \mathbf{Q}_k' \mathbf{K}_t \mathbf{A}_m \boldsymbol{\mu} \right| / R_{t,T}(\mathbf{w}), \\
&= \left| \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\mu}' \mathbf{A}_k' \mathbf{Q}_k' \mathbf{K}_t \mathbf{A}_m \boldsymbol{\mu} \right| / R_{t,T}(\mathbf{w}) \\
&= \left| \sum_{k=1}^{M_T} w^k \boldsymbol{\mu}' \mathbf{A}_k' \mathbf{Q}_k' \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu} \right| / R_{t,T}(\mathbf{w}) \\
&\leq \left[\boldsymbol{\mu}' \sum_{m=1}^{M_T} w^m \mathbf{A}_m \mathbf{Q}_m \sum_{k=1}^{M_T} w^k \mathbf{Q}_k \mathbf{A}_k \boldsymbol{\mu} \boldsymbol{\mu}' \mathbf{A}(\mathbf{w}) \mathbf{K}_t^2 \mathbf{A}(\mathbf{w}) \boldsymbol{\mu} / R_{t,T}^2(\mathbf{w}) \right]^{\frac{1}{2}} \\
&\leq \left[k_{\max} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\mu}' \mathbf{A}_k \mathbf{Q}_k \mathbf{Q}_m \mathbf{A}_m \boldsymbol{\mu}' / R_{t,T}(\mathbf{w}) \right]^{\frac{1}{2}},
\end{aligned}$$

and

$$\begin{aligned}
&\sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\mu}' \mathbf{A}_k \mathbf{Q}_k \mathbf{K}_t \mathbf{Q}_m \mathbf{A}_m \boldsymbol{\mu} / R_{t,T}(\mathbf{w}) \\
&\leq 2k_{\max} \xi_{t,T}^{-1} \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\mu}' (\mathbf{A}_k \mathbf{Q}_k \mathbf{Q}_m \mathbf{A}_m + \mathbf{A}_m \mathbf{Q}_m \mathbf{Q}_k \mathbf{A}_k) \boldsymbol{\mu} \\
&\leq k_{\max} \xi_{t,T}^{-1} \boldsymbol{\mu}' \boldsymbol{\mu} \max_{1 \leq k \leq M_T} \max_{1 \leq m \leq M_T} \zeta(\mathbf{A}_k \mathbf{Q}_k \mathbf{Q}_m \mathbf{A}_m) \\
&\leq k_{\max} \xi_{t,T}^{-1} \tilde{h}^2 \boldsymbol{\mu}' \boldsymbol{\mu} \rightarrow 0.
\end{aligned}$$

Then $\Delta_{21} = o_p(1)$ and thus (A.18) is proved.

For (A.19),

$$\begin{aligned}
\tilde{L}_{t,T}(\mathbf{w}) - \tilde{R}_{t,T}(\mathbf{w}) &= \boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - 2 \boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} - \text{tr} \left(\tilde{\mathbf{P}}(\mathbf{w})' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\Omega} \right) \\
&\equiv \Delta_{31} + \Delta_{32} + \Delta_{33}.
\end{aligned}$$

Firstly, for Δ_{32} ,

$$\begin{aligned}
&\left| \frac{\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon}}{\tilde{R}_{t,T}(\mathbf{w})} \right| \\
&\leq \left[\frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}(\mathbf{w})' \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} (\boldsymbol{\mu}' \tilde{\mathbf{A}}(\mathbf{w})' \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu})}{\tilde{R}_{t,T}^2(\mathbf{w})} \right]^{\frac{1}{2}} \\
&\leq \left[\frac{\boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon}}{\tilde{R}_{t,T}(\mathbf{w})} \right]^{\frac{1}{2}}.
\end{aligned}$$

Besides, we can obtain that

$$\begin{aligned}
& \boldsymbol{\varepsilon}' \tilde{\mathbf{P}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}(\mathbf{w}) \boldsymbol{\varepsilon} \\
& \leq \boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}' \mathbf{Q}'(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\varepsilon} \\
& \quad + \boldsymbol{\varepsilon}' \mathbf{T}'(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon} + 2|\boldsymbol{\varepsilon}' \mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon}| \\
& \quad + 2|\boldsymbol{\varepsilon}' \mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\varepsilon}| + 2|\boldsymbol{\varepsilon}' \mathbf{Q}(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon}|.
\end{aligned} \tag{A.26}$$

Terms on the right side of (A.26) follow that:

$$\begin{aligned}
2|\boldsymbol{\varepsilon}' \mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon}| & \leq [\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' \mathbf{T}'(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon}]^{\frac{1}{2}}, \\
2|\boldsymbol{\varepsilon}' \mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\varepsilon}| & \leq [\boldsymbol{\varepsilon}' \mathbf{P}'(\mathbf{w}) \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' \mathbf{Q}'(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\varepsilon}]^{\frac{1}{2}}, \\
2|\boldsymbol{\varepsilon}' \mathbf{Q}(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon}| & \leq [\boldsymbol{\varepsilon}' \mathbf{T}'(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' \mathbf{Q}(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\varepsilon}]^{\frac{1}{2}}.
\end{aligned}$$

Furthermore,

$$\begin{aligned}
\boldsymbol{\varepsilon}' \mathbf{Q}'(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\varepsilon} & = \sum_{k=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\varepsilon}' \mathbf{Q}_k \mathbf{K}_t \mathbf{Q}_m \boldsymbol{\varepsilon} \\
& \leq \sum_{m=1}^{M_T} \sum_{m=1}^{M_T} w^k w^m \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} \zeta(\mathbf{Q}_k \mathbf{K}_t \mathbf{Q}_m) \leq k_{\max} \tilde{h}^2 \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}.
\end{aligned}$$

Similarly, $\boldsymbol{\varepsilon}' \mathbf{T}'(\mathbf{w}) \mathbf{K}_t \mathbf{T}(\mathbf{w}) \boldsymbol{\varepsilon} \leq k_{\max} \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) \tilde{h}^2 \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}$. Combining the proof of Δ_{11} in (A.17), we obtain Δ_{31} and Δ_{32} are both $o_p(1)$. The proof of Δ_{33} is similar to that of (A.18) and is omitted here. Thus, (A.19) is proved.

For (A.20), given Assumption 6' and (A.15), we have

$$\begin{aligned}
& \Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\mu}' \tilde{\mathbf{A}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}}{\tilde{R}_{t,T}(\mathbf{w})} \right| > \delta \right\} \\
& \leq \Pr \left\{ \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \boldsymbol{\mu}' \tilde{\mathbf{A}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} \right| > \delta \tilde{\xi}_{t,T} \right\} \\
& \leq \sum_{m=1}^{M_T} \Pr \left\{ \left| \boldsymbol{\mu}' \tilde{\mathbf{A}}'_t(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\varepsilon} \right| > \delta \tilde{\xi}_{t,T} \right\} \\
& \leq \delta^{-2} \tilde{\xi}_{t,T}^{-2} \sum_{m=1}^{M_T} \mathbb{E} \left(\boldsymbol{\mu}' \tilde{\mathbf{A}}'_t(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\varepsilon} \right)^2 \\
& = \delta^{-2} \tilde{\xi}_{t,T}^{-2} \sum_{m=1}^{M_T} \text{tr} \left(\boldsymbol{\mu}' \tilde{\mathbf{A}}'_t(\mathbf{w}_m^0) \mathbf{K}_t \boldsymbol{\Omega} \mathbf{K}_t \tilde{\mathbf{A}}'_t(\mathbf{w}_m^0) \boldsymbol{\mu} \right) \\
& \leq k_{\max} \delta^{-2} \tilde{\xi}_{t,T}^{-2} \zeta(\boldsymbol{\Omega}) \sum_{m=1}^{M_T} \tilde{R}_T(\mathbf{w}_m^0) \rightarrow 0.
\end{aligned}$$

Thus, (A.20) is proved. And the proof of (A.21) is similar and is omitted here. Given (A.17)-(A.21), Theorem 2 is valid. ■

3 Appendix A.3

Proof of Theorem 2'. First recall that the normality assumption in the proof of Theorem 1 is to make $\Omega^{-1/2}\boldsymbol{\varepsilon}$ be a vector of independent variables. Now, we directly assume that $\boldsymbol{\varepsilon}$ is a vector of independent variables. From the proof steps in Appendix A.1, it is straightforward to know that when the normality assumption in Theorem 1 is replaced by the assumption that $\boldsymbol{\varepsilon}$ is a vector of independent variables, the conclusion of Theorem 1 still holds. In addition, when $\boldsymbol{\varepsilon}$ is a vector of independent variables, Ω is a diagonal matrix and so Assumption 7 holds automatically. Hence, to prove Theorem 2', we need only to verify Assumption 5.

It is seen that

$$\begin{aligned} \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\tilde{R}_{t,T}(\mathbf{w}) - R_{t,T}(\mathbf{w})|}{R_{t,T}(\mathbf{w})} &\leq \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\boldsymbol{\mu}' \tilde{\mathbf{A}}'(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{A}}(\mathbf{w}) \boldsymbol{\mu} - \boldsymbol{\mu}' \mathbf{A}'(\mathbf{w}) \mathbf{K}_t \mathbf{A}(\mathbf{w}) \boldsymbol{\mu}|}{R_{t,T}(\mathbf{w})} \\ &\quad + \sup_{\mathbf{w} \in \mathcal{H}_T} \frac{|\text{tr}(\tilde{\mathbf{P}}(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}'(\mathbf{w}) \Omega) - \text{tr}(\mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \Omega)|}{R_{t,T}(\mathbf{w})} \\ &\equiv \Delta_5 + \Delta_6. \end{aligned} \tag{A.27}$$

From the proof of $\Delta_{21} = o_p(1)$ in Appendix A.2, we have

$$\Delta_5 = o(1), \quad a.s.. \tag{A.28}$$

In addition,

$$\begin{aligned} &\text{tr}(\tilde{\mathbf{P}}(\mathbf{w}) \mathbf{K}_t \tilde{\mathbf{P}}'(\mathbf{w}) \Omega) - \text{tr}(\mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \Omega) \\ &= \text{tr}((\mathbf{P}(\mathbf{w}) - \mathbf{Q}(\mathbf{w}) + \mathbf{T}(\mathbf{w})) \mathbf{K}_t (\mathbf{P}(\mathbf{w}) - \mathbf{Q}(\mathbf{w}) + \mathbf{T}(\mathbf{w}))' \Omega) - \text{tr}(\mathbf{P}(\mathbf{w}) \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \Omega), \end{aligned}$$

which, along with Assumptions 6', 8 and 12, implies

$$\Delta_6 = o(1), \quad a.s.. \tag{A.29}$$

From (A.27)-(A.29), we can verify Assumption 5. This completes the proof. ■

4 Appendix A.4

Proof of Theorem 3. For TVJMA, we have $\mathbf{D} \equiv \sum_{m=1}^{M_T} w^m \mathbf{D}_m = \mathbf{Q}(\mathbf{w}) + \mathbf{I}_T$. Denote \mathbf{X}_{m^c} as a matrix consisting of columns of \mathbf{X} except for \mathbf{X}_m , i.e, $\mathbf{X}_{m^c} = \mathbf{X} \boldsymbol{\Pi}'_{m^c}$ with a selection matrix $\boldsymbol{\Pi}_{m^c}$, $\boldsymbol{\beta}_t^m = \boldsymbol{\Pi}_m \boldsymbol{\beta}_t$, and $\boldsymbol{\beta}_t^{m^c} = \boldsymbol{\Pi}_{m^c} \boldsymbol{\beta}_t$. Then, let Ξ be an $M_T \times M_T$ matrix with the (i, j) -th element:

$$\Xi_{ij} = (\boldsymbol{\varepsilon} + \mathbf{X}_{i^c} \boldsymbol{\beta}_t^{i^c})' (\mathbf{I}_T - \mathbf{P}_i) (\mathbf{Q}_i \mathbf{K}_t + \mathbf{K}_t \mathbf{Q}_j + \mathbf{Q}_i \mathbf{K}_t \mathbf{Q}_j) (\mathbf{I}_T - \mathbf{P}_j) (\boldsymbol{\varepsilon} + \mathbf{X}_{j^c} \boldsymbol{\beta}_t^{j^c}),$$

and thus

$$\begin{aligned} CV_{t,T}(\mathbf{w}) &= (\mathbf{Y} - \mathbf{P}(\mathbf{w})\mathbf{Y})'\mathbf{K}_t(\mathbf{Y} - \mathbf{P}(\mathbf{w})\mathbf{Y}) + \mathbf{w}'\Xi\mathbf{w} \\ &\equiv CM_{t,T}(\mathbf{w}) + \mathbf{w}'\Xi\mathbf{w}, \end{aligned}$$

where $CM_{t,T}(\mathbf{w}) \equiv (\mathbf{Y} - \mathbf{P}(\mathbf{w})\mathbf{Y})'\mathbf{K}_t(\mathbf{Y} - \mathbf{P}(\mathbf{w})\mathbf{Y})$. Based on (A.14)-(A.15) and Assumptions 3, 12-14, we have

$$\begin{aligned} &(\boldsymbol{\varepsilon} + \mathbf{X}_{i_c}\boldsymbol{\beta}_t^{i_c})'(\mathbf{I}_T - \mathbf{P}_i)(\mathbf{Q}_i\mathbf{K}_t + \mathbf{K}_t\mathbf{Q}_j + \mathbf{Q}_m\mathbf{K}_t\mathbf{Q}_j)(\mathbf{I}_T - \mathbf{P}_j)(\boldsymbol{\varepsilon} + \mathbf{X}_{j_c}\boldsymbol{\beta}_t^{j_c}) \\ &\leq \left\| (\boldsymbol{\varepsilon} + \mathbf{X}_{i_c}\boldsymbol{\beta}_t^{i_c}) \right\| \left\| \boldsymbol{\varepsilon} + \mathbf{X}_{j_c}\boldsymbol{\beta}_t^{j_c} \right\| \times \zeta \{ (\mathbf{I}_T - \mathbf{P}_i)(\mathbf{Q}_i\mathbf{K}_t + \mathbf{K}_t\mathbf{Q}_j + \mathbf{Q}_m\mathbf{K}_t\mathbf{Q}_j)(\mathbf{I}_T - \mathbf{P}_j) \} \\ &\leq \left\| (\boldsymbol{\varepsilon} + \mathbf{X}_{i_c}\boldsymbol{\beta}_t^{i_c}) \right\| \left\| \boldsymbol{\varepsilon} + \mathbf{X}_{j_c}\boldsymbol{\beta}_t^{j_c} \right\| k_{\max} \{ 2h^* + (h^*)^2 \} \left(1 + \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) \right)^2 \\ &= O_p(1). \end{aligned}$$

Thus, for any \mathbf{w} , $\mathbf{w}'\Xi\mathbf{w} = O_p(1)$.

Next, we will show that $\|\sqrt{Th}(\widehat{\boldsymbol{\beta}}_t(\widehat{\mathbf{w}}_t) - \boldsymbol{\beta}_t)\| = O_p(1)$. It is seen that

$$\begin{aligned} CM_{t,T}(\mathbf{w}) &= (\mathbf{Y} - \mathbf{P}(\mathbf{w})\mathbf{Y})'\mathbf{K}_t(\mathbf{Y} - \mathbf{P}(\mathbf{w})\mathbf{Y}) \\ &= (\boldsymbol{\mu} + \boldsymbol{\varepsilon} - \mathbf{X}\widehat{\boldsymbol{\beta}}_t(\mathbf{w}))'\mathbf{K}_t(\boldsymbol{\mu} + \boldsymbol{\varepsilon} - \mathbf{X}\widehat{\boldsymbol{\beta}}_t(\mathbf{w})) \\ &= \boldsymbol{\varepsilon}'\mathbf{K}_t\boldsymbol{\varepsilon} + (\widehat{\boldsymbol{\beta}}_t(\mathbf{w}) - \boldsymbol{\beta}_t)'\mathbf{X}'\mathbf{K}_t\mathbf{X}(\widehat{\boldsymbol{\beta}}_t(\mathbf{w}) - \boldsymbol{\beta}_t) - 2\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{X}(\widehat{\boldsymbol{\beta}}_t(\mathbf{w}) - \boldsymbol{\beta}_t). \end{aligned}$$

When $w^j = 1$ ($j \notin \mathcal{U}$) for any fixed time point t , we have $CM_{t,T}(\mathbf{w}) = \boldsymbol{\varepsilon}'\mathbf{K}_t\boldsymbol{\varepsilon} + \eta_{t,T}^j$, where \mathcal{U} is a set of under-fitted models and $\eta_{t,T}^j \equiv (\widehat{\boldsymbol{\beta}}_t^j - \boldsymbol{\beta}_t)'\mathbf{X}'\mathbf{K}_t\mathbf{X}(\widehat{\boldsymbol{\beta}}_t^j - \boldsymbol{\beta}_t) - 2\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{X}(\widehat{\boldsymbol{\beta}}_t^j - \boldsymbol{\beta}_t)$. With Assumptions 8, 13-14 and Proposition A.1 of Chen & Hong (2012), it is shown that

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_t^j - \boldsymbol{\beta}_t &= \boldsymbol{\Pi}_j'(\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\mathbf{Y} - \boldsymbol{\beta}_t \\ &= \boldsymbol{\Pi}_j'(\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\left(\mathbf{X}\left(\boldsymbol{\beta}_t + \frac{s-t}{T}\boldsymbol{\beta}_t^{(1)} + O(h^2)\right) + \boldsymbol{\varepsilon}\right) - \boldsymbol{\beta}_t \\ &= \boldsymbol{\Pi}_j'(\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\boldsymbol{\varepsilon} + h \int_{-1}^1 k(u)udu O_p(\zeta(\boldsymbol{\Pi}_j\mathbb{E}\mathbf{X}_t\mathbf{X}_t')) + O(h^2) \\ &= O_p(1/\sqrt{Th}) + O(h^2), \quad t \in [Th, T - Th], \end{aligned}$$

where $\boldsymbol{\beta}_t^{(1)}$ is the first derivative of $\boldsymbol{\beta}_t$. Then, we have

$$\begin{aligned} \eta_{t,T}^j &= ((\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\boldsymbol{\varepsilon} + O(h^2))'\boldsymbol{\Pi}_j\mathbf{X}'\mathbf{K}_t\mathbf{X}\boldsymbol{\Pi}_j'((\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\boldsymbol{\varepsilon} + O(h^2)) \\ &\quad - 2\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{X}\boldsymbol{\Pi}_j'((\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\boldsymbol{\varepsilon} + O(h^2)) \\ &= -\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{X}^j(\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)^{-1}\mathbf{X}^{j'}\mathbf{K}_t\boldsymbol{\varepsilon} + O_p(h^4\zeta(\mathbf{X}^{j'}\mathbf{K}_t\mathbf{X}^j)) \\ &= O_p(1) + O_p(Th^5) = O_p(1), \end{aligned}$$

based on $h = cT^{-\lambda}$ for $\frac{1}{5} \leq \lambda < 1$, where $0 < c < \infty$. If t in the interior region $[Th, T - Th]$, $\int_{-1}^1 k(u)udu = 0$ with Assumption 8. If t in the right boundary region $[T - Th, T]$,

$\int_{-1}^c k(u)udu \neq 0$ with Assumption 8, and thus we have $\widehat{\beta}_t^j - \beta_t = \Pi_j'(\mathbf{X}^{j'} \mathbf{K}_t \mathbf{X}^j)^{-1} \mathbf{X}^{j'} \mathbf{K}_t \boldsymbol{\varepsilon} + h \int_{-1}^c k(u)udu O_p(\zeta(\Pi_j \mathbb{E} \mathbf{X}_t \mathbf{X}_t')) + O(h^2) = O_p(1/\sqrt{Th}) + O(h)$, and $\eta_{t,T}^j = O_p(1) + O_p(Th^3) = O_p(1)$, which is based on $h = cT^{-\lambda}$ for $\frac{1}{3} \leq \lambda < 1$, where $0 < c < \infty$. The similar result holds for the left boundary region $[1, Th]$. Thus, $CM_{t,T}(\widehat{\mathbf{w}}_t) \leq \boldsymbol{\varepsilon}' \mathbf{K}_t \boldsymbol{\varepsilon} + \eta_{t,T}^j + O_p(1)$ for $t/T \in [0, 1]$. This implies that

$$\begin{aligned} \eta_{t,T}^j &\geq CM_{t,T}(\widehat{\mathbf{w}}_t) - \boldsymbol{\varepsilon}' \mathbf{K}_t \boldsymbol{\varepsilon} \\ &= (\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t)' \mathbf{X}' \mathbf{K}_t \mathbf{X} (\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t) - 2\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} (\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t). \end{aligned} \quad (\text{A.30})$$

From (A.30), we have

$$\begin{aligned} &\zeta_{\min}(\Psi_{t,T}) \|\sqrt{Th}(\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t)\|^2 \\ &\leq (\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t)' \mathbf{X}' \mathbf{K}_t \mathbf{X} (\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t) \\ &\leq \eta_{t,T}^j + 2\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} (\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t) \\ &\leq \eta_{t,T}^j + 2 \left\| T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} \right\| \left\| \sqrt{Th}(\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t) \right\|, \end{aligned}$$

and thus

$$\begin{aligned} &\zeta_{\min}(\Psi_{t,T}) \left[\left\| \sqrt{Th}(\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t) \right\| - \zeta_{\min}^{-1}(\Psi_{t,T}) \left\| T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} \right\| \right]^2 \\ &\leq \eta_{t,T}^j + \zeta_{\min}^{-1}(\Psi_{t,T}) \left\| T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} \right\|^2. \end{aligned} \quad (\text{A.31})$$

From (A.31), we have

$$\begin{aligned} \sqrt{Th} \|\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t\| &\leq \left\{ \zeta_{\min}^{-1}(\Psi_{t,T}) [\eta_{t,T}^j + \zeta_{\min}^{-1}(\Psi_{t,T}) \|T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X}\|^2] \right\}^{1/2} \\ &\quad + \zeta_{\min}^{-1}(\Psi_{t,T}) \|T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X}\|, \end{aligned}$$

and

$$\begin{aligned} \sqrt{Th} \|\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t\| &\geq - \left\{ \zeta_{\min}^{-1}(\Psi_{t,T}) [\eta_{t,T}^j + \zeta_{\min}^{-1}(\Psi_{t,T}) \|T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X}\|^2] \right\}^{1/2} \\ &\quad + \zeta_{\min}^{-1}(\Psi_{t,T}) \|T^{-1/2} h^{-1/2} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X}\|. \end{aligned}$$

Therefore, we have $\sqrt{Th} \|\widehat{\beta}_t(\widehat{\mathbf{w}}_t) - \beta_t\| = O_p(1)$, with Assumption 14. ■

5 Appendix A.5

Following Vogt (2012), the process $\{Y_t\}$ is locally stationary if for each rescaled time point $\tau \in [0, 1]$ there exists an associated strictly stationary process $Y_t(\tau)$ with $\|Y_t - Y_t(\tau)\| = O_p(h + \frac{1}{T})$. Thus, for every time t , we can replace \mathbf{Y}_L , in the neighbourhood of t (i.e., $[t - Th, t + Th]$), by a strictly stationary process $Y_t(\tau)$ with a small cost $O_p(h + \frac{1}{T})$, where $\tau = t/T$. We replace $Y_t(\tau)$ with x_t to simplify the notation. Denote $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ as the minimum and maximum eigenvalues of matrix A , respectively.

Before proving Theorem 4, we need to prove the following four lemmas:

Lemma 1. Under Assumptions 18 and 20, for any $q > 0$ and all $\theta > 0$,

$$\mathbb{E}\lambda_{\min}^{-q}\left(\widehat{\mathbf{R}}_{t,T}(r_1)\right) = O(r_1^{(2+\theta)q}), \quad (\text{A.32})$$

where $\widehat{\mathbf{R}}_{t,T}(r_1) = \frac{1}{(T-r_1)h} \sum_{j=r_1}^{T-1} \mathbf{x}_j^{(r_1)} k_{jt} \mathbf{x}_j^{(r_1)'}$ and $\mathbf{x}_j^{(r_1)} \equiv (x_j, \dots, x_{j-r_1+1})'$ in the augmented regression model.

Proof of Lemma 1. For notational simplicity, define $\mathbf{x}_j = \mathbf{x}_j^{(r_1)}$ and $\mathbf{A} = \begin{pmatrix} 1 & a_1 & \cdots & a_{r_1-1} \\ 0 & 1 & & \vdots \\ \vdots & & \ddots & a_1 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$.

Following the spirit of Lemma 1 in Ing and Wei (2003), we consider the following transformation of \mathbf{x}_j , $\phi_j = \mathbf{A}k_{jt}^{1/2}\mathbf{x}_j = \mathbf{B}_t\mathbf{x}_j = \boldsymbol{\varrho}_j + \boldsymbol{\varsigma}_j$ for any fixed time point t , where $\boldsymbol{\varrho}_j = (\varrho_j, \dots, \varrho_{j-r_1+1})'$, and $\boldsymbol{\varsigma}_j = (\varsigma_{j1}, \dots, \varsigma_{jr_1})'$ with ς_{ji} , $1 \leq i \leq r_1$, being a linear combination of $\boldsymbol{\varrho}_l$, $l \leq j - r_1$. It is easy to obtain the following results:

(F1) $\boldsymbol{\varrho}_j$ is independent of $\{\boldsymbol{\varsigma}_{l_1}, \boldsymbol{\varrho}_{l_2}\}$ for $l_1 \leq j$ and $l_2 \leq j - r_1$,

(F2) $\lambda_{\min}^{-1}(\sum_{j=r_1}^{T-1} \mathbf{x}_j k_{jt} \mathbf{x}_j') \leq \lambda_{\max}(\mathbf{B}_t' \mathbf{B}_t) \lambda_{\min}^{-1}(\sum_{j=r_1}^{T-1} \phi_j \phi_j')$,

(F3) $\lambda_{\max}(\mathbf{B}_t' \mathbf{B}_t) = O(1)$ for any fixed time point t .

In view of (F2) and (F3), (A.32) follows from

$$\mathbb{E} \left((T - r_1)^q h^q \lambda_{\min}^{-q} \left(\sum_{j=r_1}^{T-1} \phi_j \phi_j' \right) \right) \leq C \left(r_1^{(2+\theta)} \right)^q, \quad (\text{A.33})$$

which is the same as Eq (2.6) in Ing and Wei (2003). With (F1), the proof of (A.33) is similar to Eq (2.6) in Ing and Wei (2003). To save the space, we omit the proof of (A.33). \blacksquare

Denote $\mathbf{R}_t(r_1) = \mathbb{E}\mathbf{x}_t \mathbf{x}_t'$ and $\|\mathbf{C}\|^2 = \lambda_{\max}(\mathbf{C}'\mathbf{C})$ as the maximum eigenvalue of the matrix $\mathbf{C}'\mathbf{C}$.

Lemma 2. Under Assumptions 20 and 21, and $\sup_{-\infty < t < \infty} \mathbb{E}|\boldsymbol{\varepsilon}_t|^{2q} < \infty$ for some $q \geq 2$, we have

$$\mathbb{E} \|\widehat{\mathbf{R}}_{t,T}(r_1) - \mathbf{R}_t(r_1)\|^q \leq C \left(\frac{r_1^2}{(T - r_1)h} \right)^{q/2}, \quad (\text{A.34})$$

where C is some positive constant.

Proof of Lemma 2. Note that

$$\mathbb{E} \|\widehat{\mathbf{R}}_{t,T}(r_1) - \mathbf{R}_t(r_1)\|^q \leq \frac{r_1^q}{r_1^2} \sum_{i=1}^{r_1} \sum_{j=1}^{r_1} \mathbb{E} |\widehat{\gamma}_{i,j,t} - \gamma_{i-j,t}|^q, \quad (\text{A.35})$$

where $\widehat{\gamma}_{i,j,t}$ and $\gamma_{i-j,t}$ denote the (i, j) components of $\widehat{\mathbf{R}}_{t,T}(r_1)$ and $\mathbf{R}_t(r_1)$, respectively. With Proposition 1 in Mathematical Appendix of Chen and Hong (2012), we obtain that

$$\mathbb{E}|\widehat{\gamma}_{i,j,t} - \gamma_{i-j,t}| = \mathbb{E}\left|\sum_{s=r_1}^{T-1} x_{s-i+1}k_{st}x_{s-j+1} - \sum_{s=r_1}^{T-1} \mathbb{E}x_{s-i+1}k_{st}x_{s-j+1}\right| = O(((T-r_1)h)^{-1/2}).$$

Then $\mathbb{E}\|\widehat{\mathbf{R}}_{t,T}(r_1) - \mathbf{R}_t(r_1)\|^q \leq C \frac{r_1^q}{r_1^2} O(((T-r_1)h)^{-q/2}) = C \left(\frac{r_1^2}{(T-r_1)h}\right)^{q/2}$. This proof is completed. \blacksquare

Lemma 3. *Under Assumption 20, if $\sup_{-\infty < t < \infty} \mathbb{E}(|\boldsymbol{\varepsilon}_t|^q) < \infty$ for $q \geq 2$, then for $1 \leq p_m \leq r_1$ with $r_1 \leq T-1$,*

$$\mathbb{E}\left\|\frac{1}{\sqrt{(T-r_1)h}} \sum_{j=r_1}^{T-1} k_{jt} \mathbf{x}_j^{(m)} \boldsymbol{\varepsilon}_{j+1}\right\|^q \leq C(p_m)^{q/2},$$

where C is some positive constant, and p_m is the lag order of dependent variables in m th candidate model.

Proof of Lemma 3. Following the spirit of Eq (3.8) in Ing and Wei (2003), it is shown that

$$\mathbb{E}\left\|\frac{1}{\sqrt{T-r_1}} \sum_{j=r_1}^{T-1} k_{jt} \mathbf{x}_j^{(m)} \boldsymbol{\varepsilon}_{j+1}\right\|^q \leq p_m^{q/2} p_m^{-1} \sum_{l=0}^{k-1} \mathbb{E}\left\{((T-r_1)h)^{-\frac{q}{2}} \left|\sum_{j=r_1}^{T-1} k_{jt} \mathbf{x}_{j-l} \boldsymbol{\varepsilon}_{j+1}\right|^q\right\}.$$

With Assumption 12 and the convexity of $x^{q/2}$, $x > 0$, there exists some constant C satisfying

$$\mathbb{E}\left(\frac{1}{(T-r_1)h} \sum_{j=r_1}^{T-1} \mathbf{x}_{j-l}^2\right)^{q/2} \leq \frac{1}{T-r_1} \sum_{j=r_1}^{T-1} \mathbb{E}\{|\mathbf{x}_{j-l}|^q\} \leq C.$$

Thus, this lemma is proved. \blacksquare

Lemma 4. *If Assumptions 15, 18 and 19 hold, $r_1^{6+\delta} = O(T)$ for some $\delta > 0$, and $\sup_{-\infty < t < \infty} \mathbb{E}(|\boldsymbol{\varepsilon}_t|^{2q_1}) < \infty$ for some $q_1 \geq 2$, then for any $0 < q < q_1$, $\mathbb{E}\|\widehat{\mathbf{R}}_{t,T}^{-1}(r_1)\|^q \leq C$, and $\mathbb{E}\|\widehat{\mathbf{R}}_{t,T}^{-1}(r_1) - \mathbf{R}_t^{-1}(r_1)\|^{q/2} \leq C \left(\frac{r_1^2}{(T-r_1)h}\right)^{q/4}$, where C is some positive constant.*

Proof of Lemma 4. Given Lemma 1, we obtain that

$$\|\widehat{\mathbf{R}}_{t,T}^{-1}(r_1) - \mathbf{R}_t^{-1}(r_1)\|^q \leq \|\widehat{\mathbf{R}}_{t,T}^{-1}(r_1)\|^q \|\widehat{\mathbf{R}}_{t,T}(r_1) - \mathbf{R}_t(r_1)\|^q \|\mathbf{R}_t^{-1}(r_1)\|^q,$$

almost surely for large T . Based on the Holder's inequality and Lemma 2, we have

$$\mathbb{E}\|\widehat{\mathbf{R}}_{t,T}^{-1}(r_1) - \mathbf{R}_t^{-1}(r_1)\|^q \leq C \left(\mathbb{E}\|\widehat{\mathbf{R}}_{t,T}(r_1) - \mathbf{R}_t(r_1)\|^{q_1}\right)^{q/q_1} (r_1^{2+\theta})^q \leq C \left(\frac{r_1^{6+2\theta}}{(T-r_1)h}\right)^{q/2},$$

for sufficiently large T . Set $2\theta \leq \delta$, and with the assumption that $r_1^{6+\delta} = O(T)$, $\mathbb{E}||\widehat{\mathbf{R}}_{t,T}^{-1}(r_1)|| \leq C$ is obtained. Moreover, since the Cauchy-Schwarz inequality gives

$$\mathbb{E}||\widehat{\mathbf{R}}_{t,T}^{-1}(r_1) - \mathbf{R}_t^{-1}(r_1)||^{q/2} \leq C(\mathbb{E}||\widehat{\mathbf{R}}_{t,T}^{-1}(r_1)||^q)^{1/2}(\mathbb{E}||\widehat{\mathbf{R}}_{t,T}(r_1) - \mathbf{R}_t(r_1)||^q)^{1/2},$$

Lemma 4 is proved. ■

Proof of Theorem 4. The proof of Theorem 4 is similar to the proof of Theorem 3.1 in Zhang et al. (2013). First, substitute $V_{t,T}(\mathbf{w})$, $\widetilde{V}_{t,T}(\mathbf{w})$, $\xi_{t,T}^*$, $\widetilde{\xi}_{t,T}^*$, $\sigma^2 I_T$ and “in probability”, for $R_{t,T}(\mathbf{w})$, $\widetilde{R}_{t,T}(\mathbf{w})$, $\xi_{t,T}$, $\widetilde{\xi}_{t,T}$, $\mathbf{\Omega}$ and “a.s.”, respectively in Theorem 2 and its proof. To prove Theorem 4, we need to verify that

$$\xi_{t,T}^{*-1} \widetilde{h} \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} = o_p(1), \quad (\text{A.36})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \xi_{t,T}^{*-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}(\mathbf{w})' \mathbf{P}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} = o_p(1), \quad (\text{A.37})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \xi_{t,T}^{*-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}(\mathbf{w}) \boldsymbol{\varepsilon} = o_p(1), \quad (\text{A.38})$$

and

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\mu}' \widetilde{\mathbf{A}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}}{\widetilde{V}_{t,T}(\mathbf{w})} \right| = o_p(1). \quad (\text{A.39})$$

Since $\boldsymbol{\mu}' \mathbf{K}_t \boldsymbol{\varepsilon}$ is unrelated to \mathbf{w} , we can equally prove

$$\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \frac{\boldsymbol{\mu}' \widetilde{\mathbf{P}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}}{\widetilde{V}_{t,T}(\mathbf{w})} \right| = o_p(1), \quad (\text{A.40})$$

instead of (A.39). According to (A.19), for proving (A.40), we only need to verify

$$\xi_{t,T}^{*-1} \sup_{\mathbf{w} \in \mathcal{H}_T} \left| \boldsymbol{\mu}' \widetilde{\mathbf{P}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} \right| = o_p(1). \quad (\text{A.41})$$

Considering that $\{\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_T\}$ is i.i.d and $\mathbb{E} \boldsymbol{\varepsilon}_i^4 < \infty$, we have that $\boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} = O_p(T)$. Then given (A.16) and Assumption 16, (A.36) is proved.

Before the proof of (A.37), we verify some equations first. With Lemma 3 and Assumptions 18-19, we have

$$\begin{aligned} & T^{-1} h^{-1} \mathbb{E}(\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{Y}_L \mathbf{Y}_L' \mathbf{K}_t \boldsymbol{\varepsilon}) \\ &= T^{-1} h^{-1} \mathbb{E} \left[\boldsymbol{\varepsilon}' \mathbf{K}_t \left(\mathbf{Y}_L(\tau) + O_p\left(h + \frac{1}{T}\right) \right) \left(\mathbf{Y}_L(\tau) + O_p\left(h + \frac{1}{T}\right) \right)' \mathbf{K}_t \boldsymbol{\varepsilon} \right] \\ &= T^{-1} h^{-1} \mathbb{E}(\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{Y}_L(\tau) \mathbf{Y}_L'(\tau) \mathbf{K}_t \boldsymbol{\varepsilon}) + O\left(h + \frac{1}{T}\right) \end{aligned}$$

$$\begin{aligned}
&= O(r_1) + O\left(h + \frac{1}{T}\right) \\
&= O(r_1).
\end{aligned} \tag{A.42}$$

Uniformly for $m = 1, \dots, M_T$, (A.42) is valid in m th candidate model at any given time point t . Thus by Markov's inequality,

$$T^{-1}h^{-1}r_1^{-1}\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{Y}_L\mathbf{Y}_L'\mathbf{K}_t\boldsymbol{\varepsilon}' = O_p(1). \tag{A.43}$$

Also, by Assumption 17, we have

$$T^{-1}h^{-1}\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{X}^*\mathbf{X}^{*'}\mathbf{K}_t\boldsymbol{\varepsilon} = O_p(1). \tag{A.44}$$

Thus we have

$$T^{-1}h^{-1}\gamma^{-1}\boldsymbol{\varepsilon}'\mathbf{K}_t\mathbf{X}\mathbf{X}'\mathbf{K}_t\boldsymbol{\varepsilon} = O_p(1). \tag{A.45}$$

By Lemma 4 and Assumptions 18-19, we have

$$Th\mathbb{E}(\zeta((\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1})) = Th\mathbb{E}\left[\zeta\left(\left(\mathbf{Y}_L'(\tau)\mathbf{K}_t\mathbf{Y}_L(\tau) + O\left(h + \frac{1}{T}\right)\right)^{-1}\right)\right] = O(1). \tag{A.46}$$

By Markov's equality, we have

$$Th\zeta\left((\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1}\right) = O_p(1). \tag{A.47}$$

Let $\mathbf{J}_t = (\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1}\mathbf{Y}_L'\mathbf{K}_t\mathbf{X}^*(\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1/2}$. By Rao (1973), it can be shown that

$$\begin{aligned}
(\mathbf{X}'\mathbf{K}_t\mathbf{X})^{-1} &= \begin{pmatrix} (\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1} + \mathbf{J}_t\mathbf{J}_t' & -\mathbf{J}_t(\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1/2} \\ -(\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1/2}\mathbf{J}_t & (\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1} \end{pmatrix} \\
&= \begin{pmatrix} (\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + 2\begin{pmatrix} \mathbf{J}_t\mathbf{J}_t' & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1} \end{pmatrix} \\
&\quad - \begin{pmatrix} \mathbf{J}_t\mathbf{J}_t' & \mathbf{J}_t(\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1/2} \\ (\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1/2}\mathbf{J}_t & (\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1} \end{pmatrix}.
\end{aligned} \tag{A.48}$$

Thus

$$\begin{aligned}
&\zeta\left((\mathbf{X}'\mathbf{K}_t\mathbf{X})^{-1}\right) \leq \zeta\left((\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1}\right) + 2\max\left\{\zeta\mathbf{J}_t\mathbf{J}_t', \zeta\left((\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1}\right)\right\} \\
&\leq \zeta\left((\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1}\right) + 2\max\left\{\zeta\left((\mathbf{Y}_L'\mathbf{K}_t\mathbf{Y}_L)^{-1}\right)\zeta\left(T^{-1}(\mathbf{X}^{*'}\mathbf{K}_t\mathbf{X}^*)\right)\zeta(T^{-1}(\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1}), \right. \\
&\quad \left. T^{-1}\zeta\left(T^{-1}(\mathbf{X}^{*'}\mathbf{M}_t\mathbf{X}^*)^{-1}\right)\right\}.
\end{aligned} \tag{A.49}$$

Combining (A.47) and (A.49), we have

$$Th\zeta\left((\mathbf{X}'\mathbf{K}_t\mathbf{X})^{-1}\right) = O_p(1). \tag{A.50}$$

Getting back to the proof of (A.37), we notice that $\mathbf{K}_t \mathbf{P}' \mathbf{P} \mathbf{K}_t = \{z_{ij}\}_{1 \leq i, j \leq T}$, where

$$z_{ij} \equiv k_{it} k_{jt} \mathbf{X}_i \sum_{s=1}^T (\mathbf{X}' \mathbf{K}_s \mathbf{X})^{-1} \mathbf{X}'_s k_{is} k_{js} \mathbf{X}_s (\mathbf{X}' \mathbf{K}_s \mathbf{X})^{-1} \mathbf{X}'_j.$$

Notice that for any (i, j) such that $|i - t| > 2Th$ or $|j - t| > 2Th$, we obtain that $z_{ij} = 0$.

Since $\frac{1}{Th} \mathbf{X}' \mathbf{K}_s \mathbf{X} \xrightarrow{p} \mathbf{R}_s \equiv \mathbb{E} \mathbf{X}'_s \mathbf{X}_s$, we have

$$z_{ij} \xrightarrow{p} \frac{k_{it} k_{jt}}{T^2 h^2} \mathbf{X}_i \left(\sum_{s=1}^T \mathbf{R}_s^{-1} \mathbf{X}'_s k_{is} k_{js} \mathbf{X}_s \mathbf{R}_s^{-1} \right) \mathbf{X}'_j \equiv z_{ij}^*.$$

By Assumption 21, it is straightforward to obtain that $\mathbf{R}_s = \mathbf{R}_t + O(|t - s|/T)$ and $k_{it} = 0$ for $|i - t| > Th$. With the fact that $k_{is} = k_{si}$, we have

$$\begin{aligned} z_{ij}^* &= k_{it} k_{jt} \mathbf{X}_i \mathbf{R}_t^{-1} \frac{1}{T^2 h^2} \sum_{s=1}^T \mathbf{X}'_s k_{is} k_{js} \mathbf{X}_s \mathbf{R}_t^{-1} \mathbf{X}'_j + O_p \left(\frac{1}{T^2 h} \right) \\ &\leq k_{\max} k_{it} k_{jt} \mathbf{X}_i \mathbf{R}_t^{-1} \frac{1}{T^2 h^2} \sum_{s=1}^T \mathbf{X}'_s k_{si} \mathbf{X}_s \mathbf{R}_t^{-1} \mathbf{X}'_j \\ &\xrightarrow{p} k_{\max} k_{it} k_{jt} \mathbf{X}_i \mathbf{R}_t^{-1} \frac{1}{Th} \mathbf{R}_i \mathbf{R}_t^{-1} \mathbf{X}'_j \\ &= k_{\max} k_{it} k_{jt} \mathbf{X}_i \mathbf{R}_t^{-1} \frac{1}{Th} \mathbf{R}_t \mathbf{R}_t^{-1} \mathbf{X}'_j + O_p \left(\frac{1}{T^2 h} \right) \\ &= \frac{k_{\max}}{Th} k_{it} \mathbf{X}_i \mathbf{R}_t^{-1} \mathbf{X}'_j k_{jt}, \quad \text{if } i, j \in [t - Th, t + Th], \end{aligned}$$

and then as $T \rightarrow \infty$,

$$\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}' \mathbf{P} \mathbf{K}_t \boldsymbol{\varepsilon} \leq \frac{k_{\max}^2}{Th} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} \mathbf{R}_t^{-1} \mathbf{X}' \mathbf{K}_t \boldsymbol{\varepsilon} \leq \frac{k_{\max}^2 \zeta(\mathbf{R}_t^{-1})}{Th} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} \mathbf{X}' \mathbf{K}_t \boldsymbol{\varepsilon}.$$

For \mathbf{R}_t^{-1} , since $\frac{1}{Th} \mathbf{X}' \mathbf{K}_t \mathbf{X} \xrightarrow{p} \mathbf{R}_t$, and combining this with (A.49), we have that

$$Th \zeta(\mathbf{R}_t^{-1}) = O_p(1). \quad (\text{A.51})$$

Then given (A.45), we have

$$\frac{1}{Th \gamma} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{X} \mathbf{R}_t^{-1} \mathbf{X}' \mathbf{K}_t \boldsymbol{\varepsilon} = O_p(1). \quad (\text{A.52})$$

Therefore,

$$\gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}' \mathbf{P} \mathbf{K}_t \boldsymbol{\varepsilon} = O_p(1). \quad (\text{A.53})$$

Given Assumption 16 that $\gamma \xi_{t,T}^{*-1} = o_p(1)$, we have $\xi_{t,T}^{*-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}' \mathbf{P} \mathbf{K}_t \boldsymbol{\varepsilon} = o_p(1)$. Similarly, we can obtain

$$\gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'_m \mathbf{P}_m \mathbf{K}_t \boldsymbol{\varepsilon} = O_p(1), \quad (\text{A.54})$$

and $\xi_{t,T}^{*-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'_m \mathbf{P}_m \mathbf{K}_t \boldsymbol{\varepsilon} = o_p(1)$. Then if $m \neq j$, we have

$$\gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'_m \mathbf{P}_j \mathbf{K}_t \boldsymbol{\varepsilon} \leq \gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t (\mathbf{P}'_m \mathbf{P}_m + \mathbf{P}'_j \mathbf{P}_j) \mathbf{K}_t \boldsymbol{\varepsilon} / 2 = O_p(1), \quad (\text{A.55})$$

and $\xi_{t,T}^{*-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'_m \mathbf{P}_j \mathbf{K}_t \boldsymbol{\varepsilon} = o_p(1)$. Given (A.54) and (A.55), we have

$$\begin{aligned} \sup_{\mathbf{w} \in \mathcal{H}_T} \gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} &= \sup_{\mathbf{w} \in \mathcal{H}_T} \gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \sum_{m=1}^{M_T} \sum_{j=1}^{M_T} w^m w^j \mathbf{P}'_m \mathbf{P}_j \mathbf{K}_t \boldsymbol{\varepsilon} \\ &\leq \max_{1 \leq j \leq M_T} \max_{1 \leq m \leq M_T} \gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'_m \mathbf{P}_j \mathbf{K}_t \boldsymbol{\varepsilon} = O_p(1). \end{aligned} \quad (\text{A.56})$$

This completes the proof of (A.37). Similarly, we obtain (A.38).

For (A.41), we obtain that

$$\begin{aligned} &\sup_{\mathbf{w} \in \mathcal{H}_T} \left| \boldsymbol{\mu}' \tilde{\mathbf{A}}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} \right| \\ &\leq \sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\mu}' \mathbf{P}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}| + \sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\mu}' \mathbf{Q}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}| + \sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\mu}' \mathbf{T}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon}| \\ &= \sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \boldsymbol{\mu}| + \sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\mu}| + \sup_{\mathbf{w} \in \mathcal{H}_T} |\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{T}'(\mathbf{w}) \boldsymbol{\mu}| \\ &\leq \sup_{\mathbf{w} \in \mathcal{H}_T} (\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} \boldsymbol{\mu}' \boldsymbol{\mu})^{1/2} + \sup_{\mathbf{w} \in \mathcal{H}_T} (\boldsymbol{\varepsilon}' \mathbf{K}_t \boldsymbol{\mu}' \mathbf{Q}'(\mathbf{w}) \mathbf{K}_t \mathbf{Q}(\mathbf{w}) \boldsymbol{\mu})^{1/2} \\ &\quad + \sup_{\mathbf{w} \in \mathcal{H}_T} (\boldsymbol{\varepsilon}' \mathbf{K}_t \boldsymbol{\mu}' \mathbf{T}(\mathbf{w}) \mathbf{K}_t \mathbf{T}'(\mathbf{w}) \boldsymbol{\mu})^{1/2} \\ &\leq \sup_{\mathbf{w} \in \mathcal{H}_T} (\boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}'(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon} \boldsymbol{\mu}' \boldsymbol{\mu})^{1/2} + (k_{\max}^2 \tilde{h}^2 \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} \boldsymbol{\mu}' \boldsymbol{\mu})^{1/2} + \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) (k_{\max}^2 \tilde{h}^2 \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} \boldsymbol{\mu}' \boldsymbol{\mu})^{1/2} \\ &= \sup_{\mathbf{w} \in \mathcal{H}_T} \xi_{t,T}^* ((\gamma \xi_{t,T}^{*-2} \boldsymbol{\mu}' \boldsymbol{\mu}) \gamma^{-1} \boldsymbol{\varepsilon}' \mathbf{K}_t \mathbf{P}(\mathbf{w}) \mathbf{P}(\mathbf{w}) \mathbf{K}_t \boldsymbol{\varepsilon})^{1/2} + \left(1 + \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) \right) (k_{\max}^2 \tilde{h}^2 \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} \boldsymbol{\mu}' \boldsymbol{\mu})^{1/2} \\ &= \left(\xi_{t,T}^* ((\gamma \xi_{t,T}^{*-2} \boldsymbol{\mu}' \boldsymbol{\mu}) \times O_p(1))^{1/2} \right) \\ &\quad + \left(1 + \max_{1 \leq m \leq M_T} \zeta(\mathbf{P}_m) \right) \xi_{t,T}^* \left((T \gamma^{-1} \tilde{h}) (k_{\max}^2 \gamma \xi_{t,T}^{*-1}) (\xi_{t,T}^{*-1} \tilde{h} \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}) (T^{-1} \boldsymbol{\mu}' \boldsymbol{\mu}) \right)^{1/2}, \end{aligned} \quad (\text{A.57})$$

where the last step is based on (A.56). Given Assumption 3, Assumption 16, (A.16), (A.36), and (A.37), we obtain (A.41). This completes the proof. ■