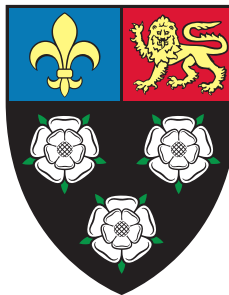


# Semiparametric Characteristics-based Models of Asset Returns



Shaoran Li

Faculty of Economics  
University of Cambridge

This dissertation is submitted for the degree of  
*Doctor of Philosophy*

King's College

April 2021



I would like to dedicate this thesis to my loving parents and my beloved wife.



## **Declaration**

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements. This dissertation contains fewer than 65,000 words including appendices, bibliography, footnotes, tables and equations and has fewer than 150 figures.

Shaoran Li  
April 2021



## Acknowledgements

I am indebted to my Ph.D. supervisor Oliver Linton for his constant encouragement, guidance and support. I have learned enormously from him. His kindness to people, patience to students and passion for researching all cheer me up. I also profit from discussions with other academic staff at the Faculty of Economics, especially Richard Smith, Alexei Onatski and Debopam Bhattacharya, for their patient help and insightful comments. Special thanks go to Sophie Jianghong Song, who makes me feel the sense of belonging at Cambridge.

I also benefit dramatically from conversations and discussions with many senior researchers. I thank Gregory Connor for his excellent ideas, patient guidance and enormous support; I thank Chaohua Dong for his priceless inspiration; I thank Xiaohong Chen for her warm words and enlightening advice during my clueless periods; I thank Liao Yuan for his valuable comments; I thank Yingying Li for her considerate help; I thank Hsiao Cheng for sharing his experience; I thank Degui Li and Jia Chen for enjoyable discussions; I thank Merrick Zhen Li, Chen Wang, Jerome Simons and Huasheng Song for their continued cheering.

I appreciate the weekly Econometrics Workshop/Seminar at the Faculty of Economics, through which this thesis is extensively improved. I thank all workshop attendees, in particular, Weiguang Liu, Ondrej Tobek and Weilun Zhou. I learn rewardingly from seminar speakers, especially Anders Kock and Joachim Freyberger. I also appreciate the prestigious King's College for its splendid dining hall and magnificent chapel, where I met and hosted many cool notables.

I would like to thank all my friends at the University of Cambridge and chauffeurs on the river of Cam, with whom the journey of my Ph.D. would be more colorful.

I am grateful to the China Scholarship Council, Cambridge Trust and Faculty of Economics for their generous financial support.

I am most grateful to my beloved family. I cannot fully express how thankful and appreciative I am to my parents and wife. I thank my parents, Defeng Li and Yanfang Bai, for all the love and trust they always give. I thank my wife Shuyi, as always, for her company.

Chapter 1 is joint work with Chaohua Dong. Chapter 2 is joint work with Shuyi Ge and Oliver Linton. Chapter 3 is joint work with Gregory Connor and Oliver Linton.



## **Abstract**

# **Semiparametric Characteristics-based Models of Asset Returns**

*Shaoran Li*

*King's College*

*University of Cambridge*

This thesis, which includes three chapters, studies asset-specific characteristics such as capitalization, book-to-market ratio etc., and their implications on assets prices and portfolio management. This thesis selects characteristics that have prediction powers on assets excess returns and specifies a flexible regression model, including linear, non-linear and pairwise interactive parts. This thesis further analyses whether characteristics are relevant as mispricing components and factor loadings in an asset pricing factor model. Finally, this thesis develops an optimal portfolio selection method based on the constructed characteristics-based asset pricing model. Methodologies in this thesis are mainly proposed for two popular questions in financial econometrics, namely, high dimensional analysis and the approximation of uni-variate and multi-variate unknown functions. The tools extended by this thesis are B-splines and orthogonal series, and multi-variate unknown functions are approximated by tensor products. In terms of high dimensional problems, which are caused by both abundant financial data and diverging B-splines bases used to approximate unknown functions, they are solved by LASSO-style selection model and power enhanced hypothesis tests. The details of the three chapters are summarized below:

## **Specification LASSO and an Application in Financial Markets**

This chapter proposes the method of Specification-LASSO in a flexible semi-parametric regression model that allows for the interactive effects between different covariates. Specification-LASSO extends LASSO and Adaptive Group LASSO to achieve both relevant variable selection and model specification. Specification-LASSO also gives preliminary estimates that facilitate the estimation of the regression model. Monte Carlo simulations show that the Specification-LASSO can accurately specify partially linear additive models with interactive effects. Finally, the proposed methods are applied in an empirical study, which examines the topic proposed by [Freyberger et al. \(2020b\)](#), arguing that firms' sizes may have interactive effects with other security-specific characteristics, which can explain the stocks excess returns together.

## **Dynamic Peer Groups of Arbitrage Characteristics**

This chapter proposes an asset pricing factor model constructed with semi-parametric characteristics-based mispricing and factor loading functions. We approximate the unknown functions by B-splines sieve where the number of B-splines coefficients is diverging. We estimate this model and test the existence of the mispricing function by a power enhanced hypothesis test. The enhanced test solves the low power problem caused by diverging B-splines coefficients, with the strengthened power approaches one asymptotically. We also investigate the structure of mispricing components through Hierarchical K-means Clusterings. We apply our methodology to CRSP (Center for Research in Security Prices) and Compustat data for the US stock market with one-year rolling windows during 1967-2017. This empirical study shows the presence of mispricing functions in certain time blocks. We also find that distinct clusters of the same characteristics lead to similar arbitrage returns, forming a "peer group" of arbitrage characteristics.

## **A Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection**

This paper develops a two-step semiparametric methodology for portfolio weight selection for characteristics-based factor-tilt and factor-timing investment strategies. We build upon the expected utility maximization framework of [Brandt \(1999\)](#) and [Ait-Sahalia and Brandt](#)

(2001). We assume that assets returns obey a characteristics-based factor model with time-varying factor risk premia as in [Ge et al. \(2020\)](#). We prove under our return-generating assumptions that an approximately optimal portfolio can be established using a two-step procedure in a market with a large number of assets. The first step finds optimal factor-mimicking sub-portfolios using a quadratic objective function over linear combinations of characteristics-based factor loadings. The second step dynamically combines these factor-mimicking sub-portfolios based on a time-varying signal, using the investor's expected utility as the objective function. We develop and implement a two-stage semiparametric estimator. We apply it to CRSP (Center for Research in Security Prices) and FRED (Federal Reserve Economic Data) data and find excellent in-sample and out-sample performance consistent with investors' risk aversion levels.



# Contents

<b>List of Figures</b>	<b>xvii</b>
<b>List of Tables</b>	<b>xix</b>
<b>1 Specification LASSO and an Application in Financial Markets</b>	<b>1</b>
1.1 Introduction . . . . .	4
1.2 Model Setup . . . . .	9
1.3 Methodology . . . . .	11
1.3.1 Variables and Model Selection by Specification-LASSO . . . . .	11
1.3.2 Estimation . . . . .	16
1.4 Theoretical results . . . . .	16
1.4.1 Assumption . . . . .	17
1.5 Simulation study . . . . .	19
1.6 Empirical Study . . . . .	21
1.6.1 Introduction . . . . .	21
1.6.2 Data Description . . . . .	22
1.6.3 Variable Selection and Model Specification . . . . .	22
1.6.4 Selection Results . . . . .	24
1.7 Conclusion . . . . .	26
<b>2 Dynamic Peer Groups of Arbitrage Characteristics</b>	<b>29</b>
2.1 Introduction . . . . .	32
2.2 Model Setup . . . . .	36
2.3 Estimation . . . . .	37
2.3.1 B-splines Approximation . . . . .	37
2.3.2 Two-Step Projected-PCA . . . . .	39
2.4 Power-enhanced Tests . . . . .	41

2.5	Hierarchical K-Means Clustering . . . . .	44
2.6	Asymptotic Properties . . . . .	48
2.6.1	Consistency Assumptions . . . . .	48
2.6.2	Main Results . . . . .	49
2.7	Numerical Study . . . . .	50
2.7.1	Data Generation . . . . .	50
2.7.2	Model Misspecification . . . . .	51
2.7.3	Robustness Under Stronger Noise . . . . .	55
2.7.4	Number of Factors . . . . .	55
2.8	Empirical Study . . . . .	59
2.8.1	Data . . . . .	59
2.8.2	Estimation . . . . .	60
2.8.3	Power-enhanced Hypothesis Tests . . . . .	60
2.8.4	Test Results . . . . .	61
2.8.5	Dynamic Peer Groups of Arbitrage Characteristics . . . . .	64
2.9	Conclusion . . . . .	67
<b>3</b>	<b>A Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection</b>	<b>69</b>
3.1	Introduction . . . . .	72
3.2	The Model of Asset Returns . . . . .	75
3.3	A Two-Step Version of the Portfolio Choice Problem . . . . .	76
3.3.1	Utility Function of the Investor . . . . .	77
3.3.2	Step 1: Factor-mimicking Sub-portfolios . . . . .	77
3.3.3	Step 2: Factor-timing Portfolio Based on Dynamic Signals . . . . .	78
3.4	Econometric Methods . . . . .	80
3.4.1	Assumptions . . . . .	80
3.4.2	Estimation of Characteristics-based Factor Loadings . . . . .	81
3.4.3	Estimation of the First Step–Factor-tilt . . . . .	81
3.4.4	Estimation of the Second Step–Factor-timing . . . . .	83
3.4.5	Hypothesis Tests . . . . .	86
3.4.6	Theoretical Results . . . . .	86
3.5	Empirical Study . . . . .	87
3.5.1	Data Description . . . . .	87
3.5.2	In-sample Factor-mimicking Portfolios . . . . .	88
3.5.3	Utility Function . . . . .	89
3.5.4	Selection of Truncation Number . . . . .	89

---

3.5.5	Estimation of Dynamic Signals . . . . .	92
3.5.6	In-sample Performance of Factor-timing Portfolios . . . . .	93
3.5.7	Out-sample Performance of Factor-timing Portfolios . . . . .	94
3.5.8	Restrictive Sub-portfolios Weights . . . . .	99
3.6	Conclusion . . . . .	104
<b>Bibliography</b>		<b>105</b>
<b>Appendix A Proofs of Theorems</b>		<b>111</b>
A.1	Proofs . . . . .	111
A.1.1	Proofs of Chapter 1 . . . . .	111
A.1.2	Proofs of Chapter 2 . . . . .	114
A.1.3	Proofs of Chapter 3 . . . . .	120
<b>Appendix B Tables and Figures</b>		<b>123</b>
B.1	Characteristic Description . . . . .	123
B.2	Figures and Tables . . . . .	127
B.2.1	Figures of Chapter 2 . . . . .	127
B.2.2	Tables of Chapter 3 . . . . .	131





# List of Figures

3.1	The Plot of DS . . . . .	90
3.2	The Plot of TS . . . . .	91
3.3	The Plot of Trend . . . . .	91
3.4	The Plot of Ln(DY%) . . . . .	91
3.5	The Plot of RF . . . . .	92
3.6	The Plot of sub-portfolios Weights . . . . .	97
3.7	The Plot of Sub-portfolio Weights . . . . .	101
3.8	The Plot of of FF 3 factors Weights . . . . .	103
B.1	Mispricing Characteristic Curve of standardized $r_{12-2}$ and $r_{12-7}$ . . . . .	127
B.2	Mispricing Characteristic Curve of standardized PCM . . . . .	127
B.3	Mispricing Characteristic Curve of standardized ROA in 1988-1989 . . . . .	128
B.4	Mispricing Characteristic Curve of standardized BEME . . . . .	128
B.5	Mispricing Characteristic Curve of standardized LME . . . . .	128
B.6	Mispricing Characteristic Curve of standardized AT . . . . .	129
B.7	Mispricing Characteristic Curve of standardized LEV in 2002-2003 . . . . .	129
B.8	Mispricing Characteristic Curve of standardized IPM in 2004-2005 . . . . .	129
B.9	Mispricing Characteristic Curve of standardized DelGmSale in 2015-2016 . . . . .	130
B.10	Mispricing Characteristic Curve of standardized C2D in 2016-2017 . . . . .	130
B.11	Clustering of PCM 1986-1987 . . . . .	130
B.12	Clustering of IPM 2004-2005 . . . . .	131



# List of Tables

1.1	Simulation Example of S-LASSO . . . . .	20
1.2	Summary of Linear Effects of Characteristics on Assets Excess Returns . . .	24
1.3	Summary of nonlinear Effects of Characteristics on Assets Excess Returns .	25
1.4	Summary of Interactive Effects of Characteristics with Size on Assets Excess Returns . . . . .	26
2.1	Simulation Results 1 Part1 . . . . .	53
2.2	Simulation Results 1 Part2 . . . . .	54
2.3	Simulation Results 2 Part1 . . . . .	57
2.4	Simulation Results 2 Part2 . . . . .	58
2.5	Empirical Study Results . . . . .	63
2.6	First layer 1986-1987 (clusterings of $\ddot{y}_{it}$ ) . . . . .	64
2.7	Second layer 1986-1987 (clusterings of characteristic PCM) . . . . .	64
2.8	Second layer 1986-1987 (clusterings of characteristic PCM) . . . . .	65
2.9	First layer 2004-2005 (clusterings of $\ddot{y}_{it}$ ) . . . . .	66
2.10	Second layer 2004-2005 (clusterings of characteristic IPM) . . . . .	66
2.11	Second layer 2004-2005 (clusterings of characteristic IPM) . . . . .	66
3.1	Index Variable Summary . . . . .	90
3.2	Tests Summary . . . . .	90
3.3	Factor-mimicking Portfolios Summary . . . . .	90
3.4	Index Variable Summary . . . . .	92
3.5	Average Annual In-sample Results . . . . .	96
3.6	Monthly Out-sample Results Comparison . . . . .	98
3.7	Average Annual In-sample Results . . . . .	100
3.8	Average Annual In-sample Results . . . . .	102
B.1	Characteristic Details . . . . .	123

B.2	Annual Correlation Between Subportfolios and Risk Factors 1-20 . . . . .	132
B.3	Annual Correlation Between Subportfolios and Risk Factors 21-40 . . . . .	132
B.4	Annual Correlation Between Subportfolios and Risk Factors 41-50 . . . . .	132

# **Chapter 1**

## **Specification LASSO and an Application in Financial Markets**



## Abstract

This paper proposes the method of Specification-LASSO in a flexible semi-parametric regression model that allows for the interactive effects between different covariates. Specification-LASSO extends LASSO and Adaptive Group LASSO to achieve both relevant variable selection and model specification. Specification-LASSO also gives preliminary estimates that facilitate the estimation of the regression model. Monte Carlo simulations show that the Specification-LASSO can accurately specify partially linear additive models with interactive regressors. Finally, the proposed methods are applied in an empirical study, which examines the topic proposed by [Freyberger et al. \(2020a\)](#), which argues that firms' sizes may have interactive effects with other security-specific characteristics, which can explain the stocks excess returns together.

KEYWORDS: Variable Selection; Model Selection; Interaction;

JEL CLASSIFICATION: C14; G12.

## 1.1 Introduction

In a data-rich era, researchers are more likely to suffer both "variable selection" and "specification" challenges. "Variable selection" problem is incurred due to the ease of data attainability, so vast of data are available when researchers intend to model. This seems to be trivial if the number of observations  $n$  is relatively large compared with the number of potential covariates  $P$ . However, in recent empirical studies that have large  $P$  and small  $n$ , which causes the classical analysis tool failing to work. Therefore, it is crucial to determine which subset of candidate variables should be considered. Meanwhile, another challenge comes from the model specification, as one may be dazzled to choose a suitable model from a model zoo. In general, all parametric analyses have the risk of misspecification. Thus, nonparametric analysis is introduced to relax the functional form restrictions. Although this helps to increase the model flexibility, the "curse of dimensionality" causes the extremely low convergence rate of estimation when the dimension of independent variables is more than three.

Suppose we observe a sample of data  $\{(Y_i, \mathbf{P}_i) : 1 \leq i \leq n\}$ , where  $i$  represents the  $i^{\text{th}}$  individual.  $\mathbf{P}_i$  is a  $P \times 1$  large dimensional vector of potential covariates where only the  $Q \times 1$  subset  $\mathbf{Q}_i$  contains relevant regressors to explain or predict the variation of  $Y_i$ , which presents a sparse model if  $Q \ll P$ .

We suppose:

$$E(Y_i | \mathbf{P}_i) = \theta_i + h(\mathbf{Q}_i), \quad i = 1, 2, \dots, n, \quad (1.1)$$

where  $\theta_i$  is the intercept whereas  $h(\mathbf{Q}_i)$  is an unknown multi-variate function of  $\mathbf{Q}_i$ . Most researchers specify an additive semi-parametric structure on  $h(\mathbf{Q}_i)$  as:

$$h(\mathbf{Q}_i) = \sum_{q=1}^Q f_q(X_{iq}), \quad (1.2)$$

where  $f_q(X_{iq})$  is an unknown uni-variate function. Models like [Equation 1.2](#) are called additive nonparametric regressions and are widely discussed by [Hastie and Tibshirani \(1990\)](#), [Linton \(1997\)](#), [Linton \(2000\)](#), and [Linton and Härdle \(1996\)](#).

The [Equation 1.2](#) avoids the curse of dimensionality by imposing an additive structure, but can be inefficient if some of the relevant covariates only have linear effects as the rate of convergence for nonparametric function  $f_q(X_{qi})$  is slower than  $O(n^{-1/2})$ .



Therefore, a partially linear additive semi-parametric model is proposed to take advantages of linear effects as:

$$h(\mathbf{Q}_i) = \theta + \sum_{l=1}^L \beta_l X_{il} + \sum_{q=L+1}^Q f_q(X_{iq}), \quad (1.3)$$

where we distinguish  $L$  linear effects from  $\mathbf{Q}_i$ , and the coefficients of linear part can be estimated at the rate of convergence  $O(n^{-1/2})$ , as discussed in Wang et al. (2007) and Ma and Yang (2011). Similar models of Equation 1.3 are also studied by Li (2000), Fan and Li (2003) and Liang et al. (2008).

Unfortunately, both additive models omit potential interactions between covariates. Pairwise interactions between covariates are quite common in both economic and financial studies.

**Example 1.1.1.** In macroeconomics, most production functions specify a interactive term of capital and labour inputs such as:

$$\text{Cobb-Douglas: } Y = \Gamma X_C^\alpha X_L^\beta + \varepsilon$$

**Example 1.1.2.** In microeconomics, Deaton and Muellbauer (1980) document the utility model of a household ( $Y$ ) containing interactions between eating and drinking ( $X_E, X_D$ ) for foodstuffs, housing and fuel ( $X_H, X_F$ ) for shelters, and television and sports ( $X_T, X_S$ ) for entertainment.

$$Y = m_{ED}(X_E, X_D) + m_{HF}(X_H, X_F) + m_{TS}(X_T, X_S) + \varepsilon$$

**Example 1.1.3.** In environment studies, Dong et al. (2019) study effects of  $CO_2$  and solar irradiance (SI) on the global sea level ( $Y_{SL}$ ) rise. They specify the model as:

$$Y_{SL} = m(X_{CO_2}, X_{SI}) + \varepsilon,$$

and they verify the interactive effects between  $CO_2$  ( $X_{CO_2}$ ) and solar irradiance ( $X_{SI}$ ) through empirical results.

**Example 1.1.4.** In finance, Freyberger et al. (2020a) argue that assets returns at time  $t$  is predictable by stock characteristics, such as capitalization and book-to-market ratio, at  $t - 1$

as

$$Y_t = \boldsymbol{\theta}_t + \underbrace{\sum_{q \neq s}^Q m_{qs}(\mathbf{X}_{qt-1} \cdot \mathbf{X}_{st-1})}_{\text{interaction with firm size Xs}} + \underbrace{\sum_{q=1}^Q m_q(\mathbf{X}_{qt-1})}_{\text{uni-variate}},$$

and they find significant effects of interactions between firms' sizes and other characteristics. In this paper, we will revisit this study using our methods.

Interactions among covariates refer to the circumstance that marginal effects of the  $j^{\text{th}}$  variable  $X_j$  on  $Y$  are determined by other relevant covariates. [Sperlich et al. \(2002\)](#) illustrate the importance of interactions in the additive model, and propose a marginal integration style estimation and test methods to solve the potential interactions in the model. However, their methods cannot be applied to a high-dimensional case, not only due to the enormous workload but also the failure of estimation when  $P > n$ . From the above examples and [Sperlich et al. \(2002\)](#), we can conclude that higher-order interactions are barely discussed due to both the curse of dimensionality and interpretation issues. In this paper, we mainly discuss pairwise interactions among variables, although our methods can be easily extended to higher-order interactions.

Based on the aforementioned research and examples, it is more reasonable to expand  $h(\mathbf{Q}_i)$  in [Equation 1.2](#) to three components, including linear, nonlinear and pairwise interactive parts.

Compared with specifying the structure of an unknown multivariate function  $h(\mathbf{Q}_i)$ , selecting relevant variables under a high-dimensional setting is more widely discussed. The most popular way for achieving this goal is LASSO (Least Absolute Shrinkage and Selection Operator) style variables selection methods. [Tibshirani \(1996\)](#) proposes this method to perform both variable selection and regularization in the linear model under high-dimensional cases.

$$\min_{\boldsymbol{\alpha}} \sum_{i=1}^n \left( Y_i - \sum_{j=1}^W \alpha_j X_{ij} \right)^2 + \lambda_n \sum_{j=1}^W |\alpha_j|, \quad (1.4)$$

In [Equation 1.4](#),  $\lambda_n$  is a data driving tuning parameter, and the attractive property of LASSO is that it can achieve initial selection by shrinking some  $\alpha = 0$  and estimation even if  $P \gg n$ . A necessary condition for consistent selection of LASSO is discussed by [Zhao and Yu \(2006\)](#) and [Zou \(2006\)](#), which is called irrepresentable condition (discussed in [subsection 1.4.1](#)).

This condition restricts the correlation between relevant and irrelevant components to be relatively small.

To relax this condition, [Zou \(2006\)](#) proposes Adaptive LASSO, which can achieve consistent selection under mild conditions:

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n \left( Y_i - \sum_{j=1}^W \beta_j X_{ij} \right)^2 + \lambda_n \sum_{j=1}^W \hat{w}_j |\beta_j|, \quad (1.5)$$

where the weight  $\hat{w}_j$  is data-dependent and typically chosen as  $\hat{w}_j = |\hat{\alpha}_j|^{-\gamma}$  for some  $\gamma > 0$ , and  $\hat{\alpha}_j$  is a preliminary consistent estimate in [Equation 1.4](#).  $X_j$  with a smaller estimate  $\hat{\alpha}_j$  will be penalized more severely, and the variable with  $\alpha = 0$  will be smoothed out.

As for selecting nonparametric functions, [Lin and Zhang \(2006\)](#) introduce COSSO (Component Selection and Smoothing Operator), where they consider the model selection in a general setting of the smoothing spline analysis of variance (SS-ANOVA) framework, shown as:

$$h(\mathbf{X}_i) = b + \sum_{j=1}^d f_j(X_i^{(j)}) + \sum_{j < k} f_{jk}(X_i^{(j)}, X_i^{(k)}) + \dots$$

This model can provide large flexibility in terms of the form of nonparametric functions, such as higher dimensional functions. However, in COSSO, it only works under  $P < n$ , which means variables considered are not allowed to exceed the number of observations. Furthermore, [Lin and Zhang \(2006\)](#) do not give a detailed discussion of the selection of the linear part. Finally, this selection model is not facilitated with the initial estimation. All of these issues will be solved by our method.

Moreover, [Huang et al. \(2010a\)](#) introduce a selection and estimation method of an additive nonparametric model inspired by both group LASSO as in [Yuan and Lin \(2006\)](#) and adaptive group LASSO as in [Wang and Leng \(2008\)](#). They use a linear combination of B-splines basis  $\phi_k, 1 \geq k \geq m_n$  to approximate any potential unknown function as:

$$f_{nj}(x) = \sum_{k=1}^{m_n} \beta_{jk} \phi_k(x).$$

Next, they consider the penalized least squares criterion

$$L_n(\mu, \boldsymbol{\beta}_n) = \sum_{i=1}^n [Y_i - \mu - \sum_{j=1}^P \sum_{k=1}^{m_n} \beta_{jk} \phi_k(X_{ij})]^2 + \lambda_n \sum_{j=1}^P \hat{w}_{nj} \|\boldsymbol{\beta}_{nj}\|_2,$$

where  $\lambda_n$  is a tuning parameter while  $\|\beta_{nj}\|_2$  is the  $L^2$  norm of the  $j^{\text{th}}$  coefficient vector  $\beta_{nj} = (\beta_{j1}, \dots, \beta_{jm_n})^\top$ , and

$$\hat{w}_{nj} = \begin{cases} \|\tilde{\beta}_{nj}\|_2^{-1} & \text{for } \|\tilde{\beta}_{nj}\|_2 > 0 \\ \infty & \text{for } \|\tilde{\beta}_{nj}\|_2 = 0 \end{cases},$$

where  $\tilde{\beta}_{nj}$  is an initial and consistent estimate. [Huang et al. \(2010a\)](#) also compare adaptive group LASSO model with COSSO by [Lin and Zhang \(2006\)](#), concluding that, when the number of observations is small, adaptive group LASSO has much higher accuracy in terms of selecting relevant variables in the semi-parametric additive model.

This paper proposes a Specification-LASSO (S-LASSO) for both variables selection and model specification of a **partially linear additive semi-parametric model with interactions**, which can be applied when  $P > n$ . S-LASSO can achieve variables selection, model specification, and initial estimation at the same time. S-LASSO firstly use levels, B-splines bases and pairwise tensor products of all potentially relevant variables to approximate linear, nonlinear and interactive effects, respectively, and then it extends a two-step procedure to give consistent selection.

In the first step, S-LASSO uses ordinary LASSO to consider all bases indifferently to attain the initial selection and estimates. In the second step, S-LASSO clusters these bases into different groups according to linear, nonlinear and interactive parts, and then an adaptive group LASSO is applied to give a final selection and estimation results. The estimates from the first step help the second step to set group-specific penalty-weighting parameters, which leads to the consistency of selection.

In the empirical work, we employ S-LASSO to study a characteristics-based asset pricing model. In [Freyberger et al. \(2020a\)](#), they assume assets excess returns can be predicted by security-relevant characteristics and their interaction with the firm's size:

$$\mathbf{Y}_t = \boldsymbol{\theta}_t + \underbrace{\sum_{q \neq s}^Q m_{qs}(\mathbf{X}_{qt-1} \cdot \mathbf{X}_{st-1})}_{\text{interaction with firm size Xs}} + \underbrace{\sum_{q=1}^Q m_q(\mathbf{X}_{qt-1})}_{\text{uni-variate}},$$

where  $\mathbf{Y}_t$  is a  $n \times 1$  vector of assets excess returns at time  $t$  while  $\mathbf{X}_{jt-1}$  is a  $n \times 1$  vector of asset-specific characteristic at time  $t - 1$ . However, they fail to consider the potential linear effects of characteristics, which have a quicker convergence rate and less computational burden. Furthermore, they analyse interactive effects by specifying the form of pairwise

interaction as  $\mathbf{X}_{qt-1} \cdot \mathbf{X}_{st-1}$  (elementwise product), which is quite restrictive since  $m_{qs}(\mathbf{X}_{qt-1} \cdot \mathbf{X}_{st-1}) \neq m_{qs}(\mathbf{X}_{qt-1}, \mathbf{X}_{st-1})$  generally. S-LASSO can overcome this limitation by considering the linear effect and not restricting the form of interactions. We will illustrate these through both simulation and empirical studies.

The rest of the paper is organized as follows. Section 2 presents the model that S-LASSO is working on; Section 3 provides procedures for S-LASSO to work; Section 3 illustrates the theoretical results; Section 4 gives simulated experiments; Section 5 demonstrates an empirical study; Section 6 concludes the paper. All proofs and other materials are arranged in the Appendix.

## 1.2 Model Setup

Suppose we observe a sample data  $(\mathbf{Y}, \mathbf{P})$ , where  $\mathbf{Y}$  presents the  $n \times 1$  vector of dependent variables while  $\mathbf{P}$  denotes the  $n \times P$  matrix of potential covariates  $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_P)$ , allowing for  $P > n$ .

We assume there is an  $n \times Q$  matrix  $\mathbf{Q} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_Q)$  that is relevant to explain or predict the variation of  $\mathbf{Y}$  and  $\mathbf{Q} \subset \mathbf{P}$ . We restrict that  $Q$  is fixed, whereas  $P$  is diverging as sample size  $n \rightarrow \infty$ . We propose a sparse structure by assuming  $Q$  is relatively small as:

$$\begin{aligned} \mathbf{Y} &= \boldsymbol{\theta} + h(\mathbf{Q}) + \mathbf{U}, \\ E(\mathbf{Y}|\mathbf{P}) &= \boldsymbol{\theta} + h(\mathbf{Q}), \end{aligned} \quad (1.6)$$

where  $\mathbf{U}$  is an  $n \times 1$  vector of idiosyncratic errors  $\varepsilon_i$  with  $E(\mathbf{U}|\mathbf{P}) = \mathbf{0}$ ;  $h(\mathbf{Q})$  is a multi-variate unknown function.

We also specify a partially linear additive semi-parametric model with interactive terms on  $h(\mathbf{Q})$  as:

$$E(\mathbf{Y}|\mathbf{P}) = \boldsymbol{\theta} + h(\mathbf{Q}) = \boldsymbol{\theta} + \overbrace{\sum_{1 \leq s < s' \leq S}^S m_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'})}^{\text{interactive}} + \overbrace{\sum_{q=1}^Q m_q(\mathbf{X}_q)}^{\text{uni-variate}} \quad (1.7)$$

$$= \boldsymbol{\theta} + \overbrace{\sum_{1 \leq s < s' \leq S}^S m_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'})}^{\text{interactive}} + \overbrace{\sum_{r=1}^R m_r(\mathbf{X}_r)}^{\text{nonlinear}} + \overbrace{\sum_{l=1}^L \beta_l \mathbf{X}_l}^{\text{linear}}, \quad (1.8)$$

where  $\mathbf{X}_j$  denotes the vector of the  $j^{\text{th}}$  covariate.  $L$ ,  $R$  and  $S$  are cardinal numbers of three sets corresponding to linear effects variables, non-linear effects variables and interactive variables, respectively, which will be estimated later. The complement of  $\mathbf{X}$  that does not appear in Equation 1.7 are regarded as irrelevant variables, which should be smoothed out.

Here we have  $Q$  relevant variables in total and  $S$  of them have interactive effects with  $S \leq Q$ . Similarly,  $R$  of them have uni-variate effects with  $R \leq Q$ . Finally,  $L$  out of  $Q$  covariates have linear effects, namely,  $R + L \leq Q$ , which means we may have some covariates having only interactive effects with others.  $s$  and  $s'$  ( $s < s'$ ) is the  $s^{\text{th}}$  pair of relevant covariates that has interaction.

Meanwhile,  $m_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'})$  is an unknown bivariate nonparametric function of the  $s^{\text{th}}$  pair of relevant variables;  $m_r(\mathbf{X}_r)$  is an uni-variate unknown function of the  $r^{\text{th}}$  relevant variable;  $\beta_l$  is the coefficient of the  $l^{\text{th}}$  relevant variable.

Furthermore, we define variable sets as follows:

$$\mathcal{L} = \{\mathbf{X}_l \in \mathbf{Q} : \mathbf{X}_l \text{ has linear effects on } \mathbf{Y}\},$$

$$\mathcal{R} = \{\mathbf{X}_r \in \mathbf{Q} : \mathbf{X}_r \text{ has nonlinear effects on } \mathbf{Y}\}$$

$$\mathcal{S} = \{\mathbf{X}_s, \mathbf{X}_{s'} \in \mathbf{Q} : \mathbf{X}_s, \mathbf{X}_{s'} \text{ have interactive effects on } \mathbf{Y}\}.$$

The cardinality for each set are:  $|\mathcal{L}| = L$ ,  $|\mathcal{R}| = L$  and  $|\mathcal{S}| = S$ . Each set above is unknown to researchers and can be empty.

Equation 1.7 avoids the curse of dimensionality with fewer restrictions. Compared with conventional additive models where components are uni-variate, we allow potential covariates to interact with each other to provide more information and flexibility. We also allow for a linear part since it has a better convergence rate and less computational burden. Therefore, practitioners do not bother employing nonparametric techniques when simpler parametric methods work. The decomposition in Equation 1.7 gives considerable adaptability to mitigate possible model misspecification. We do not include higher-order interactions among covariates, but our methods can be extended accordingly.

Based on the model above, our methodology focuses on:

1. Selecting the relevant variables subset  $\mathbf{Q}$  from  $\mathbf{P}$ ;
2. Specifying the form of decomposition in Equation 1.7;

3. Giving initial estimates of Equation 1.7.

## 1.3 Methodology

This section provides the detailed procedures to select relevant variables, decompose and estimate of  $h(\mathbf{X})$ .

### 1.3.1 Variables and Model Selection by Specification-LASSO

Without external knowledge and other information, it is hard for us to determine relevant variables and the form of Equation 1.7. Therefore, all forms of entire covariates and their interactive effects should be considered, and then, a proper variable selection model can be applied to filter all possibilities. After analyzing selection results, one can examine whether the function form of each covariate is linear or not, and whether some of them have interactive effects.

To develop our methods and theoretical results, we introduce some notations and definitions. First, we illustrate spline spaces.

Similar to Schumaker (1981) and Huang et al. (2010a), we suppose that the  $j^{th}$  potential covariates  $\mathbf{X}_j$ , where  $\mathbf{X}_j$  is a  $n \times 1$  vector taking values in  $[a, b]$  as:

$$\mathbf{X}_j = (X_{1j}, X_{2j}, \dots, X_{nj})^\top, \quad \mathbf{X}_j \in \mathbf{P} \text{ and } j = 1, 2, \dots, P.$$

Furthermore,  $a, b$  are finite with  $a < b$ . Let  $\mathbf{K} = \{a = \underbrace{\kappa_0 = \kappa_0 = \dots = \kappa_0}_g < \kappa_1 < \kappa_2 < \dots < \kappa_{k_n} < \underbrace{\kappa = \kappa = \dots = \kappa}_g = b\}$  be a sequence of knots partitioning the interval  $[a, b]$  into subintervals, where  $k_n = [n^\nu]$  with  $0 < \nu < 0.5$  being a positive integer whereas  $g$  is the order of bases used. Let  $K_n = k_n + g$ , which denotes the total number of bases. For the  $i^{th}$  individual of  $\mathbf{X}_j$ , where  $j = 1, 2, \dots, P$  and  $i = 1, 2, \dots, n$ , a set of B-splines can be built in the  $L^2$  space  $\Omega_n[\mathbf{K}]$  as  $\Phi_{\mathbf{K}}(\mathbf{X}_j) = \{\phi_1(\mathbf{X}_j), \phi_2(\mathbf{X}_j), \dots, \phi_{K_n}(\mathbf{X}_j)\}$ . Next, we define a B-splines matrix:

$$\Phi_{\mathbf{K}}(\mathbf{X}_j) = \begin{pmatrix} \phi_1(X_{j1}) & \phi_2(X_{j1}) & \dots & \phi_{K_n}(X_{j1}) \\ \phi_1(X_{j2}) & \phi_2(X_{j2}) & \dots & \phi_{K_n}(X_{j2}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(X_{jn}) & \phi_2(X_{jn}) & \dots & \phi_{K_n}(X_{jn}) \end{pmatrix},$$

**Definition 1.3.1.** Define spline space  $\mathcal{H}_{g,\mathbf{K}}$  as linear combination of B-splines by:

$$\mathcal{H}_{g,\mathbf{K}} = \text{span}\{\phi_{\mathbf{K},k}, 1 \leq k \leq K_n\} = \left\{ \sum_{k=1}^{K_n} \beta_k \phi_{\mathbf{K},k} \mid \beta_k \in \mathbb{R} \text{ for } 1 \leq k \leq K_n \right\},$$

where  $g$  is the degree of those bases and  $\mathbf{K}$  is the knots sequence, and  $\beta_k$  is the  $k^{\text{th}}$  B-spline coefficient. To simplify the notation without causing confusion, we drop the sequence subscript  $\mathbf{K}$  henceforth.

Accordingly, the  $r^{\text{th}}$  unknown uni-variate function can be approximated as:

$$m_r(\mathbf{X}_r) = \Phi(\mathbf{X}_r)\boldsymbol{\beta}_r + \boldsymbol{\xi}_r,$$

where  $\boldsymbol{\beta}_r = (\beta_{r1}, \beta_{r2}, \dots, \beta_{rK_n})^\top$ , and  $\boldsymbol{\xi}_r$  is the approximation error.

Similar to spline space  $\mathcal{H}_{g,\mathbf{K}}$ , we construct another spline space  $\mathcal{D}_{g,\mathbf{D}}$  using knot sequence  $\mathbf{D}$  in interval  $[a', b']$ .

**Definition 1.3.2.** Define the **tensor product** of spline spaces  $\mathcal{H}_{g,\mathbf{K}} \otimes \mathcal{D}_{g,\mathbf{D}}$  as the family of all functions of the form:

$$f(\mathbf{x}_p, \mathbf{x}_{p'}) = \sum_{k=1}^{K_n} \sum_{d=1}^{D_n} \beta_{kd} \phi_k(\mathbf{x}_p) \mu_d(\mathbf{x}_{p'}), \text{ where } 1 < 2 < \dots < p < p' < \dots < P$$

where coefficients  $\beta_{kd}$  can be any real numbers.

Accordingly, for any covariates  $\mathbf{X}_a, \mathbf{X}_b \in \mathbf{P}$ , their potential interactive effects can be approximated as:

$$m_{ab}(\mathbf{X}_a, \mathbf{X}_b) = \sum_{k=1}^{K_n} \sum_{d=1}^{D_n} \beta_{abkd} \phi_k(\mathbf{X}_a) \mu_d(\mathbf{X}_b) + \boldsymbol{\xi}_{ab}, \quad 1 \leq a < b \leq P,$$

where  $\boldsymbol{\xi}_{ab}$  is the approximation error.

Equivalently, let

$$\Phi_{\mathbf{K}}(X_{ia}) = (\phi_1(X_{ia}), \phi_2(X_{ia}), \dots, \phi_{K_n}(X_{ia}))^\top,$$

$$\mu_{\mathbf{D}}(X_{ib}) = (\mu_1(X_{ib}), \mu_2(X_{ib}), \dots, \mu_{D_n}(X_{ib}))^\top.$$



Equivalently:

$$\begin{aligned} \Phi_{\mathbf{K}}(X_{ia}) \otimes \mu_{\mathbf{D}}(X_{ib}) &= \text{Vec} \left( \begin{bmatrix} \phi_1(X_{ia})\mu_1(X_{ib}) & \phi_1(X_{ia})\mu_2(X_{ib}) & \dots & \phi_1(X_{ia})\mu_{D_n}(X_{ib}) \\ \phi_2(X_{ia})\mu_1(X_{ib}) & \phi_2(X_{ia})\mu_2(X_{ib}) & \dots & \phi_2(X_{ia})\mu_{D_n}(X_{ib}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{K_n}(X_{ia})\mu_1(X_{ib}) & \phi_{K_n}(X_{ia})\mu_2(X_{ib}) & \dots & \phi_{K_n}(X_{ia})\mu_{D_n}(X_{ib}) \end{bmatrix} \right)^\top \\ &= (\phi_1(X_{ia})\mu_1(X_{ib}), \phi_1(X_{ia})\mu_2(X_{ib}), \dots, \phi_1(X_{ia})\mu_{D_n}(X_{ib}), \dots, \phi_{K_n}(X_{ia})\mu_{D_n}(X_{ib})). \end{aligned}$$

Then:

$$\Phi_{\mathbf{K}}(\mathbf{X}_a) \otimes \mu_{\mathbf{D}}(\mathbf{X}_b) = \begin{bmatrix} \Phi_{\mathbf{K}}(X_{1a}) \otimes \mu_{\mathbf{D}}(X_{1b}) \\ \Phi_{\mathbf{K}}(X_{2a}) \otimes \mu_{\mathbf{D}}(X_{2b}) \\ \vdots \\ \Phi_{\mathbf{K}}(X_{na}) \otimes \mu_{\mathbf{D}}(X_{nb}) \end{bmatrix}.$$

To simplify the notation without causing any confusion, we drop the sequence subscript  $\mathbf{K}$  and  $\mathbf{D}$  henceforth.

We also write tensor product coefficients as vector  $\boldsymbol{\beta}_{ab}$  as:

$$\boldsymbol{\beta}_{ab} = (\beta_{ab11}, \beta_{ab12}, \dots, \beta_{ab1D_n}, \dots, \beta_{abK_n1}, \beta_{abK_n2}, \dots, \beta_{abK_nD_n})^\top$$

The true model can be approximated as:

$$\mathbf{Y} = \boldsymbol{\theta} + \sum_{\mathbf{X}_l \in \mathcal{L}} \beta_l \mathbf{X}_l + \sum_{\mathbf{X}_r \in \mathcal{R}} \boldsymbol{\Phi}(\mathbf{X}_r) \boldsymbol{\beta}_r + \sum_{\mathbf{X}_s, \mathbf{X}_{s'} \in \mathcal{S}} \boldsymbol{\Phi}(\mathbf{X}_s) \otimes \boldsymbol{\mu}(\mathbf{X}_b) \boldsymbol{\beta}_{ss'} + \boldsymbol{\varepsilon}_n + \mathbf{U},$$

where  $\boldsymbol{\varepsilon}_n$  is the approximation error and  $\mathbf{U}$  is the  $n \times 1$  vector of idiosyncratic error  $\varepsilon_i$ .

Those non-zero coefficients are:

$$\boldsymbol{\beta}_{\mathcal{L}} = (\beta_1, \dots, \beta_L)^\top,$$

$$\boldsymbol{\beta}_{\mathcal{R}} = (\boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_R^\top)^\top,$$

$$\boldsymbol{\beta}_{\mathcal{S}} = (\boldsymbol{\beta}_{11'}^\top, \dots, \boldsymbol{\beta}_{SS'}^\top)^\top.$$

We define a non-zero coefficient vector:

$$\boldsymbol{\beta}_{P_1} = (\boldsymbol{\beta}_{\mathcal{L}}^\top, \boldsymbol{\beta}_{\mathcal{R}}^\top, \boldsymbol{\beta}_{\mathcal{S}}^\top)^\top.$$

Let  $\dim(\boldsymbol{\beta}_{P_1}) = P_1$ , where  $\dim(\cdot)$  means the dimension of any vector. We also define B-spline bases of relevant covariates as:

$$\mathbf{X}_{\mathcal{L}} = (\mathbf{X}_1, \dots, \mathbf{X}_l, \dots, \mathbf{X}_L), \quad \mathbf{X}_l \in \mathcal{L}.$$

$$\mathbf{N}(\mathbf{X}_{\mathcal{R}}) = (\boldsymbol{\Phi}(\mathbf{X}_1), \dots, \boldsymbol{\Phi}(\mathbf{X}_r), \dots, \boldsymbol{\Phi}(\mathbf{X}_R)), \quad \mathbf{X}_r \in \mathcal{R}.$$

$$\mathbf{I}(\mathbf{X}_{\mathcal{S}}) = (\boldsymbol{\Phi}(\mathbf{X}_1) \otimes \boldsymbol{\mu}(\mathbf{X}_{1'}), \dots, \boldsymbol{\Phi}(\mathbf{X}_S) \otimes \boldsymbol{\mu}(\mathbf{X}_{S'})), \quad \mathbf{X}_s \text{ and } \mathbf{X}_{s'} \in \mathcal{S}.$$

Recall that for individual  $i$ , we observe  $P$  potential covariates denoted as a vector  $\mathbf{P}_i$ , and there are  $Q$  relevant variables denoted as  $\mathbf{Q}_i$ ,  $Q \leq P$ . There are two steps for the S-LASSO to work to select  $\mathbf{Q}_i$  out of  $\mathbf{P}_i$  and to specify the model as in [Equation 1.7](#).

In the next step, our job is to put all possible linear, nonlinear and interactive forms of all potential covariates in a selection model. S-LASSO can achieve at least three goals, namely, to select all the relevant variables, to specify the model and to obtain the preliminary estimates.

*Step 1.* Substitute all possible forms of each variable and pairwise interactive terms in  $\mathbf{P}$  into LASSO selection:

$$\begin{aligned} \min_{\boldsymbol{\beta}_l, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ab}} \quad & \left\| \mathbf{Y} - \boldsymbol{\theta} - \sum_{l=1}^P \beta_l \mathbf{X}_l - \sum_{r=1}^P \boldsymbol{\Phi}(\mathbf{X}_r) \boldsymbol{\beta}_r - \sum_{a=1}^{P-1} \sum_{b>a}^P \boldsymbol{\Phi}(\mathbf{X}_a) \otimes \boldsymbol{\mu}(\mathbf{X}_b) \boldsymbol{\beta}_{ab} \right\|_2^2 \\ & + \lambda_n \left( \sum_{l=1}^P |\beta_l| + \sum_{r=1}^P \|\boldsymbol{\beta}_r\| + \sum_{a=1}^{P-1} \sum_{b>a}^P \|\boldsymbol{\beta}_{ab}\| \right) \end{aligned}$$

where  $|\beta|$  and  $\|\boldsymbol{\beta}_n\|$  are  $l_1$  norms and  $\|\boldsymbol{\beta}\|_2 \equiv (\sum_{n=1}^N |\beta_n|^2)^{1/2}$  denotes the  $l_2$  norm of any  $n \times 1$  vector  $\boldsymbol{\beta}$ .  $\lambda_n > 0$  is a data driven tuning parameter. This step provides us with preliminary information after the initial selection. However, one drawback of LASSO process is that it may leave plenty of small but non-zero coefficients. Nonetheless, the first step provide crucial hints which are helpful for discriminatory penalty in the next step.

*Step 2.* Use step 1 estimates to construct penalty weighting coefficients and substitute all bases into adaptive group LASSO:

$$\hat{\omega}_l = \begin{cases} \sqrt{N_{\mathcal{L}}} |\tilde{\beta}_l|^{-1}, & \text{if } |\tilde{\beta}_l| > 0 \\ \infty, & \text{if } \tilde{\beta}_l = 0. \end{cases}$$

$$\hat{\omega}_r = \begin{cases} \sqrt{N_{\mathcal{R}}} \|\tilde{\boldsymbol{\beta}}_r\|_2^{-1}, & \text{if } \|\tilde{\boldsymbol{\beta}}_r\|_2 > 0 \\ \infty, & \text{if } \|\tilde{\boldsymbol{\beta}}_r\|_2 = 0. \end{cases}$$

$$\hat{\omega}_{ab} = \begin{cases} \sqrt{N_{\mathcal{S}}} \|\tilde{\boldsymbol{\beta}}_{ab}\|_2^{-1}, & \text{if } \|\tilde{\boldsymbol{\beta}}_{ab}\|_2 > 0 \\ \infty, & \text{if } \|\tilde{\boldsymbol{\beta}}_{ab}\|_2 = 0. \end{cases}$$

$N_{\mathcal{L}} = L$ ,  $N_{\mathcal{R}} = R \times K_n$  and  $N_{\mathcal{S}} = \frac{S(S-1)}{2} \times (K_n)^2$  are the number of coefficients within each group as our group sizes are significantly different. We use group cardinality to control the strength of the penalty.

To eliminate the noise from step 1, we consider the adaptive group LASSO which can select variables in a group manner.

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}_l, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ab}) = \|\mathbf{Y} - \boldsymbol{\theta} - \sum_{l=1}^P \beta_l \mathbf{X}_l - \sum_{r=1}^P \boldsymbol{\Phi}(\mathbf{X}_r) \boldsymbol{\beta}_r - \sum_{a=1}^{P-1} \sum_{b>a}^P \boldsymbol{\Phi}(\mathbf{X}_a) \otimes \boldsymbol{\mu}(\mathbf{X}_b) \boldsymbol{\beta}_{ab}\|_2$$

$$+ \tilde{\lambda}_n \left( \sum_{l=1}^P \hat{\omega}_l |\beta_l| + \sum_{r=1}^P \hat{\omega}_r \|\boldsymbol{\beta}_r\|_2 + \sum_{a=1}^{P-1} \sum_{b>a}^P \hat{\omega}_{ab} \|\boldsymbol{\beta}_{ab}\|_2 \right),$$

Let  $0 \times \infty = 0$ , so groups deleted by LASSO are not selected by adaptive group LASSO for sure.  $\tilde{\lambda}_n > 0$  is a data driven tuning parameter.

After the selection by step 2, all non-zero coefficients of linear approximation are represented as  $\hat{\boldsymbol{\beta}}_{\mathcal{L}}$ ; non-zero coefficients of the approximate of nonlinear effects are shown as  $\hat{\boldsymbol{\beta}}_{\mathcal{R}}$ ; non-zero coefficients of tensor products are written as  $\hat{\boldsymbol{\beta}}_{\mathcal{S}}$ . At the same time, all the irrelevant variables or bases are smoothed out since their coefficients are zeros. Additionally, the non-zero  $\beta$ s of Step 2 is a vector  $\hat{\boldsymbol{\beta}}_{P_1}$ ,

$$\hat{\boldsymbol{\beta}}_{P_1} = (\hat{\boldsymbol{\beta}}_{\mathcal{L}}^{\top}, \hat{\boldsymbol{\beta}}_{\mathcal{R}}^{\top}, \hat{\boldsymbol{\beta}}_{\mathcal{S}}^{\top})^{\top},$$

where  $\hat{\boldsymbol{\beta}}_{\mathcal{L}} = (\hat{\beta}_1, \dots, \hat{\beta}_L)^{\top}$ ,  $\hat{\boldsymbol{\beta}}_{\mathcal{R}} = (\hat{\boldsymbol{\beta}}_1^{\top}, \dots, \hat{\boldsymbol{\beta}}_R^{\top})^{\top}$ , and  $\hat{\boldsymbol{\beta}}_{\mathcal{S}} = (\hat{\boldsymbol{\beta}}_{11'}^{\top}, \dots, \hat{\boldsymbol{\beta}}_{\hat{S}\hat{S}'}^{\top})^{\top}$ .

The model specification we obtained is:

$$h(\mathbf{Q}) = \sum_{\mathbf{X}_l \in \hat{\mathcal{L}}} \beta_l \mathbf{X}_l + \sum_{\mathbf{X}_r \in \hat{\mathcal{R}}} m_r(\mathbf{X}_r) + \sum_{\mathbf{X}_s, \mathbf{X}_{s'} \in \hat{\mathcal{S}}} m_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'}),$$

where,

$$\hat{\mathcal{L}} = \{\mathbf{X}_l \in \mathbf{Q} : |\hat{\beta}_l| > 0\},$$

$$\hat{\mathcal{R}} = \{\mathbf{X}_r \in \mathbf{Q} : \|\hat{\boldsymbol{\beta}}_r\|_2 > 0\},$$

$$\hat{\mathcal{S}} = \{\mathbf{X}_s, \mathbf{X}_{s'} \in \mathcal{Q} : \|\hat{\boldsymbol{\beta}}_{ss'}\|_2 > 0\},$$

Accordingly,  $|\hat{\mathcal{L}}| = \hat{L}$ ,  $|\hat{\mathcal{R}}| = \hat{R}$  and  $|\hat{\mathcal{S}}| = \hat{S}$ . In practice, we include covariates that are selected by both linear and nonlinear parts in the nonlinear set only since this can simplify the model further. The classification above is for theoretical proof purposes.

Next, nonlinear and interactive components are approximated by:

$$\begin{aligned} \hat{m}_r(\mathbf{X}_r) &= \boldsymbol{\Phi}(\mathbf{X}_r) \hat{\boldsymbol{\beta}}_r, 1 \leq r \leq \hat{R} \\ \hat{m}_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'}) &= \boldsymbol{\Phi}(\mathbf{X}_s) \otimes \boldsymbol{\mu}(\mathbf{X}_{s'}) \hat{\boldsymbol{\beta}}_{ss'}, 1 \leq s < s' \leq \hat{S}. \end{aligned}$$

Meanwhile, we define the matrix of irrelevant components, which are smoothed out by S-LASSO as:

$$\mathbf{X}_{\mathcal{L}^c} = (\mathbf{X}_1, \dots, \mathbf{X}_l, \dots, \mathbf{X}_{LC}), \quad \mathbf{X}_l \in \mathbf{P} \text{ but } \mathbf{X}_l \notin \mathcal{L}.$$

$$\mathbf{N}(\mathbf{X}_{\mathcal{R}^c}) = (\boldsymbol{\Phi}(\mathbf{X}_1), \dots, \boldsymbol{\Phi}(\mathbf{X}_r), \dots, \boldsymbol{\Phi}(\mathbf{X}_{RC})), \quad \mathbf{X}_r \in \mathbf{P} \text{ but } \mathbf{X}_r \notin \mathcal{R}.$$

$$\mathbf{I}(\mathbf{X}_{\mathcal{S}^c}) = (\boldsymbol{\Phi}(\mathbf{X}_1) \otimes \boldsymbol{\mu}(\mathbf{X}_{1'}), \dots, \boldsymbol{\Phi}(\mathbf{X}_{SC}) \otimes \boldsymbol{\mu}(\mathbf{X}_{SC'})), \quad \mathbf{X}_s \text{ and } \mathbf{X}_{s'} \in \mathbf{P} \text{ but } \mathbf{X}_s \text{ and } \mathbf{X}_{s'} \notin \mathcal{S}.$$

Let  $n \times P_1$  matrix  $\mathbf{Z}_1 = (\mathbf{X}_{\mathcal{L}}, \mathbf{N}(\mathbf{X}_{\mathcal{R}}), \mathbf{I}(\mathbf{X}_{\mathcal{S}}))$  represent all the relevant components and let  $\boldsymbol{\beta}_{P_1}$  be the  $P_1 \times 1$  coefficient vector of matrix  $\mathbf{Z}_1$ . Meanwhile, let  $n \times P_2$  matrix  $\mathbf{Z}_2 = (\mathbf{X}_{\mathcal{L}^c}, \mathbf{N}(\mathbf{X}_{\mathcal{R}^c}), \mathbf{I}(\mathbf{X}_{\mathcal{S}^c}))$ , denotes all the irrelevant components. Similarly, let  $\boldsymbol{\beta}_{P_2} = (\boldsymbol{\beta}_{\mathcal{L}^c}^\top, \boldsymbol{\beta}_{\mathcal{R}^c}^\top, \boldsymbol{\beta}_{\mathcal{S}^c}^\top)^\top$  be the  $P_2 \times 1$  coefficient vector of matrix  $\mathbf{Z}_2$ .

### 1.3.2 Estimation

OLS can be applied to obtain estimates:

$$\hat{\boldsymbol{\beta}}_{P_1} = (\mathbf{Z}_1^\top \mathbf{Z}_1)^{-1} \mathbf{Z}_1^\top \mathbf{Y}.$$

And

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{P_2} &= \mathbf{0}, \\ \hat{\boldsymbol{\beta}}_{P_Z} &= (\hat{\boldsymbol{\beta}}_{P_1}^\top, \hat{\boldsymbol{\beta}}_{P_2}^\top)^\top. \end{aligned}$$

## 1.4 Theoretical results

Firstly, we list some assumptions to facilitate our theoretical analysis.

### 1.4.1 Assumption

**Assumption 1.4.1.** *The noise  $\varepsilon_i$  are independent and identically distributed with  $E\varepsilon_i = 0$  and  $\text{Var}(\varepsilon) = \sigma^2$ . Furthermore, it has finite  $2k^{\text{th}}$  moment with  $E(\varepsilon_i^{2k}) < \infty$  for  $k = 1, 2, \dots, K$ .*

**Assumption 1.4.2.** *Let*

$$\mathbf{V} = \frac{1}{n}(\mathbf{Z}_1, \mathbf{Z}_2)^\top(\mathbf{Z}_1, \mathbf{Z}_2) = \begin{Bmatrix} \mathbf{V}_{Z_1Z_1} & \mathbf{V}_{Z_1Z_2} \\ \mathbf{V}_{Z_2Z_1} & \mathbf{V}_{Z_2Z_2} \end{Bmatrix}$$

*be the covariance matrix of all the components in step 1. There exist constants  $c_1, c_2, c_3$ , and  $c_4$  with  $0 \leq c_1 < c_2 \leq 1$  and  $c_3, c_4 > 0$ , such that*

$$P_1 = O(n^{c_1}), \quad (1.9)$$

$$n^{\frac{1-c_2}{2}} \min\{|\beta_l|, \|\boldsymbol{\beta}_r\|_2, \|\boldsymbol{\beta}_{ss'}\|_2\} \geq c_4, \text{ for } \beta_l, \boldsymbol{\beta}_r, \boldsymbol{\beta}_{ss'} \in \boldsymbol{\beta}_{P_1}. \quad (1.10)$$

$$P_2 = O(n^{(c_2-c_1)k}), \quad (1.11)$$

$$\lambda_{\min}(\mathbf{V}_{Z_1Z_1}) > c_3, \quad (1.12)$$

**Equation 1.9** and **Equation 1.11** control the maximum dimensions of relevant and irrelevant components respectively. **Equation 1.12** ensures that the minimum eigenvalue of relevant components matrix  $\mathbf{Z}_1$  is away from 0 to be invertible, where  $\lambda_{\min}(\mathbf{V}_{Z_1Z_1})$  indicates the smallest eigenvalue of covariance matrix  $\mathbf{V}_{Z_1Z_1}$ . Finally, **Equation 1.10** limits the decay rate of elements in  $\boldsymbol{\beta}_{P_1}$ .

**Assumption 1.4.3.**  $E(m_r(\mathbf{X}_r)) = 0$ ,  $E(m_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'})) = 0$ , given  $\mathbf{X}_j \in \mathcal{R} \cup \mathcal{S}$ .

*This assumption is for unique identification purpose.*

**Assumption 1.4.4.** *0-th, first and second derivatives of  $m_r(\mathbf{X}_r)$  and  $m_{ss'}(\mathbf{X}_s, \mathbf{X}_{s'})$  are continuous, for  $X_r \in \mathcal{R}$  and  $X_s, X_{s'} \in \mathcal{S}$ .*

This assumption is for approximation accuracy of B-splines bases and their tensor products.

**Definition 1.4.1.** Let  $\hat{\beta}$  be an estimate of  $\beta$ . Then,  $\hat{\beta}$  is **Sign Consistent** with  $\beta$ , shown as  $\hat{\beta} =_s \beta$ , if and only if

$$\text{sign}(\hat{\beta}) = \text{sign}(\beta),$$

where  $\text{sign}(\hat{\beta}) = 1$ , if  $\hat{\beta} > 0$ ;  $\text{sign}(\hat{\beta}) = -1$ , if  $\hat{\beta} < 0$ ; and  $\text{sign}(\hat{\beta}) = 0$ , if  $\hat{\beta} = 0$ . Similarly, Let  $\hat{\boldsymbol{\beta}}$  be a vector of estimates of  $\boldsymbol{\beta}$ . Then  $\hat{\boldsymbol{\beta}}$  is **Sign Consistent** with  $\boldsymbol{\beta}$ , written as  $\hat{\boldsymbol{\beta}} =_s \boldsymbol{\beta}$  if and only if each entry is **Sign Consistent**.

**Definition 1.4.2.** Let  $\hat{\beta}$  be an estimate of  $\beta$ . Then,  $\hat{\beta}$  is **Norm Consistent** with  $\beta$ , shown as  $\hat{\beta} =_0 \beta$ , if and only if

$$\text{sign}_0(\hat{\beta}) = \text{sign}_0(\beta),$$

where  $\text{sign}_0(\hat{\beta}) = 1$ , if  $\hat{\beta} \neq 0$ ;  $\text{sign}_0(\hat{\beta}) = 0$ , if  $\hat{\beta} = 0$ . Similarly, Let  $\hat{\boldsymbol{\beta}}$  be a vector of estimates of  $\boldsymbol{\beta}$ . Then  $\hat{\boldsymbol{\beta}}$  is **Norm Consistent** with  $\boldsymbol{\beta}$ , written as  $\hat{\boldsymbol{\beta}} =_0 \boldsymbol{\beta}$  if and only if each entry is **Norm Consistent**.

**Condition 1.4.1.** Let covariance matrix  $\mathbf{V}$  satisfies strong irrerepresentable condition documented by [Zhao and Yu \(2006\)](#), stating that there exists a positive constant  $P_1 \times 1$  vector  $\boldsymbol{\eta}$ , and

$$|\mathbf{V}_{Z_2 Z_1} (\mathbf{V}_{Z_1 Z_1})^{-1} \text{sign}(\boldsymbol{\beta}_{P_1})| \leq \mathbf{1} - \boldsymbol{\eta},$$

which is true element-wise.

**Condition 1.4.2.** Similarly, covariance matrix  $\mathbf{V}$  satisfies weak irrerepresentable condition , if

$$|\mathbf{V}_{Z_2 Z_1} (\mathbf{V}_{Z_1 Z_1})^{-1} \text{sign}(\boldsymbol{\beta}_{P_1})| < \mathbf{1},$$

which is true element-wise.

**Theorem 1.4.1.** Under Assumptions [1.4.1-1.4.4](#) and Condition [1.4.2](#), and let  $P_Z = P_1 + P_2$ ,  $\boldsymbol{\beta}_{P_Z} = (\boldsymbol{\beta}_{P_1}^\top, \boldsymbol{\beta}_{P_2}^\top)^\top$ , for  $\forall \lambda_n$  satisfying  $\frac{\lambda_n}{\sqrt{n}} = o(n^{\frac{c_2 - c_1}{2}})$  and  $\frac{1}{P_Z} (\frac{\lambda_n}{\sqrt{n}})^{2k} \rightarrow \infty$  for  $k = 1, 2, 3, \dots$ , then the first step of S-LASSO is sign consistent with:

$$P(\hat{\boldsymbol{\beta}}_{P_Z} =_s \boldsymbol{\beta}_{P_Z}) \geq 1 - O\left(\frac{P_Z n^k}{\lambda_n^{2k}}\right) \rightarrow 1, \text{ as } n \rightarrow \infty.$$

**Theorem 1.4.2.** Given the well-chosen number of internal knots  $k_n = \lceil n^\nu \rceil$  and under Assumptions [1.4.1-1.4.4](#) and Theorem [1.4.1](#), S-LASSO is consistent on selection relevant covariates and specification of the correct model:

$$P(\lim_{n \rightarrow \infty} \hat{\boldsymbol{\beta}}_{\mathcal{L}} =_0 \boldsymbol{\beta}_{\mathcal{L}}) \rightarrow 0,$$

$$P(\lim_{n \rightarrow \infty} \hat{\boldsymbol{\beta}}_{\mathcal{R}} =_0 \boldsymbol{\beta}_{\mathcal{R}}) \rightarrow 0,$$

$$P(\lim_{n \rightarrow \infty} \hat{\boldsymbol{\beta}}_{\mathcal{J}} =_0 \boldsymbol{\beta}_{\mathcal{J}}) \rightarrow 0.$$

## 1.5 Simulation study

We generate our model as:

$$y_i = \beta x_{1i} + m_1(x_{2i}) + m_2(x_{3i}, x_{4i}) + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\beta x_{1i} = x_{1i}$ ,  $m_1(x_{2i}) = x_{2i}^2$ ,  $m_2(x_{3i}, x_{4i}) = \sin(x_{3i} + x_{4i})$ . All three above functions are rescaled to be zero mean and unit variance. Furthermore, we generate  $P$  candidate variables  $x_{pi}$ . We have all independent variables generated from  $Uniform[-2, 2]$  and  $\varepsilon \sim N(0, \sigma^2)$ , where there are no correlations between all potential variables. Two different  $\mathbf{P}$  dimensions and three different sample sizes are tested, namely,  $P = 30, 50$  with  $n = 100, 300, 500$ .

In [Table 1.1](#), we compare the results of S-LASSO and the methods of selecting interactive effects between stock characteristics in [Freyberger et al. \(2020a\)](#) (named FNW).

We choose four evenly distributed knots to construct B-splines approximation of nonlinear effects while choosing two evenly distributed knots for each covariate to construct tensor products to keep the group size comparable. Meanwhile, for the FNW methods, we choose all the knots sequences for  $x_j$  and  $x_j \times x_{j'}$  to be 4, which are also evenly distributed, to approximate both nonlinear and interaction effects among potential covariates. The tuning parameter  $\lambda_n$ s are chosen through BIC for both steps. Here, we define BIC as:

$$BIC = n * \log(MSE) + df * \log(n),$$

where  $n$  is the number of observation and  $df$  represents the degree of freedom in LASSO procedures discussed in [Leng et al. \(2006\)](#).

Furthermore, we define the signal to noise ratio as  $R_\sigma = sd(m(\cdot))/sd(\varepsilon)$  to illustrate the robustness of S-LASSO under different noise level.

Table 1.1 Simulation Example of S-LASSO

		$\sigma=0.25$			$\sigma=0.333$			$\sigma=0.5$			
		INC	CS	MSE	INC	CS	MSE	INC	CS	MSE	
P=30	n=100	S-LASSO	52.8 (0.5)	49.4 (0.5)	0.77 (0.28)	48.2 (0.5)	44.2 (0.5)	0.83 (0.29)	32.8 (0.47)	28.4 (0.45)	0.94 (0.32)
		FNW	0 (0)	0 (0)	1.13 (0.4)	0 (0)	0 (0)	1.2 (0.41)	0.2 (0.04)	0.2 (0.04)	1.42 (0.49)
	n=300	S-LASSO	95.6 (0.21)	95 (0.22)	0.59 (0.17)	94.8 (0.22)	94 (0.24)	0.66 (0.2)	95 (0.218)	95 (0.218)	0.81 (0.2)
		FNW	0 (0)	0 (0)	0.97 (0.18)	0 (0)	0 (0)	1.02 (0.17)	0 (0)	0 (0)	1.16 (0.18)
	n=500	S-LASSO	99.8 (0.04)	99.6 (0.06)	0.53 (0.08)	98.2 (0.133)	97.2 (0.165)	0.59 (0.11)	98.4 (0.126)	98.4 (0.126)	0.75 (0.138)
		FNW	0 (0)	0 (0)	0.97 (0.14)	0 (0)	0 (0)	0.96 (0.14)	0 (0)	0 (0)	1.15 (0.15)
P=50	n=100	S-LASSO	34.8 (0.48)	31.6 (0.47)	0.88 (0.31)	34.8 (0.48)	32 (0.47)	0.89 (0.32)	24.8 (0.43)	22 (0.41)	1.01 (0.35)
		FNW	0 (0)	0 (0)	1.24 (0.48)	0 (0)	0 (0)	1.31 (0.45)	0 (0)	0 (0)	1.5 (0.41)
	n=300	S-LASSO	92.8 (0.26)	92.4 (0.27)	0.66 (0.22)	88 (0.33)	88 (0.33)	0.75 (0.26)	85.8 (0.35)	85.8 (0.35)	0.92 (0.27)
		FNW	0 (0)	0 (0)	1.01 (0.2)	0 (0)	0 (0)	1.06 (0.19)	0 (0)	0 (0)	1.22 (0.19)
	n=500	S-LASSO	98.6 (0.12)	98.4 (0.13)	0.57 (0.12)	96.8 (0.18)	96.8 (0.18)	0.64 (0.17)	96 (0.2)	96 (0.2)	0.81 (0.19)
		FNW	0 (0)	0 (0)	0.98 (0.14)	0 (0)	0 (0)	1.02 (0.16)	0 (0)	0 (0)	1.18 (0.16)

Note: This table compares the performance of S-LASSO and the method used in FNW (2020) under different sample size,  $n=100, 300, 500$ ; different number of irrelevant variables,  $P=30, 50$ ; and different levels of noise,  $R_\sigma = 4, 3, 2$ . INC represents the percentage that all the relevant covariates are correctly included in the model. CS shows the percentage of the whole model that is correctly specified, which means the model not only selects all relevant variables but also gives them a precise specification. MSE indicates the average mean squared error of all repetitions under each method. Simulations are repeated 500 times for each setting. Standard deviations are given in the parentheses.

From the results in [Table 1.1](#), S-LASSO overperforms FNW under all scenarios. Because in FNW, they treat interaction term  $x_j \times x_{j'}$  as a new variable and construct B-spline space based on this covariate. Therefore, only certain forms of pairwise interactions with input  $x_j \times x_{j'}$  can be detected. Hence, for nearly all the simulation settings, FNW can neither include all the relevant covariates nor specify the model correctly, given the interactive



function form  $\sin(x_{3i} + x_{4i})$ . However, S-LASSO employs tensor products of B-splines to approximate potential interactions and has decent accuracy on both including all relevant covariates and choosing the correct model. We use this simulation to show the limitation of FNW and demonstrate that tensor products can accommodate more comprehensive forms of interactions. Additionally, although prediction is not the primary goal of S-LASSO, it has much smaller MSE compared with FNW.

Furthermore, S-LASSO works better for small  $P$  large  $n$  circumstances, and the highest percentage of selecting the relevant covariates and specifying the correct model can be 99.8% and 99.6% individually. For the most challenging condition under  $P=50$ ,  $n=100$  and  $\sigma = 0.5$ , S-LASSO also has acceptable performance with an accuracy of 24.8% and 22% respectively.

S-LASSO is also robust under different levels of noise for all settings. As shown from different rows of [Table 1.1](#), the accuracy is similar across three different noise levels.

## 1.6 Empirical Study

### 1.6.1 Introduction

In this section, we revisit the question proposed by [Freyberger et al. \(2020a\)](#), where they try to detect the influence of firms' characteristics on stock returns non-parametrically. They specify assets returns as additive non-parametric functions of lagged corresponding assets characteristics such as book-to-market ratio, profitability, etc. Their model is:

$$E(\mathbf{Y}_{t+1}|\mathbf{W}_t) = \boldsymbol{\theta}_t + \sum_{r=1}^R m_r(\mathbf{X}_{rt}), \quad (1.13)$$

where  $\mathbf{Y}_{t+1}$  is a  $n \times 1$  vector of stock excess returns at time  $t$ ;  $\mathbf{W}_t$  is a  $n \times W$  matrix of  $W$  asset-relevant characteristics that are observed at time  $t$ . At the right hand side of [Equation 1.13](#), they select  $R$  additive non-parametric uni-variate unknown functions of characteristics that are relevant to predict stock excess returns, and  $\boldsymbol{\theta}_t$  is the intercept.

To further investigate the interactive effects between assets sizes with other characteristics, they propose a model to accommodate pairwise interactions:

$$E(\mathbf{Y}_{t+1}|\mathbf{W}_t) = \boldsymbol{\theta}_t + \sum_{r=1}^R m_r(\mathbf{X}_{rt}) + \sum_s^S m_s(\mathbf{X}_{st} \cdot \mathbf{X}_{size,t}), \quad (1.14)$$

where they consider the unknown function form taking input as  $\mathbf{X}_s \cdot \mathbf{X}_{size}$ . As discussed in Introduction and exemplified in Simulation,  $m_s(\mathbf{X}_s \cdot \mathbf{X}_{size}) \neq m_s(\mathbf{X}_s, \mathbf{X}_{size})$ , this specification of interactions may restrict the form of interactive effects to be multiplicity only. Furthermore, they do not include linear parts, which have both computational simplicity and quicker rate of convergence.

In this section, we apply S-LASSO to short rolling window data to revisit the effects of assets characteristics on stock returns and their interactive effects with firms sizes. We further divide uni-variate effects to be linear or nonlinear. The model is specified as:

$$E(\mathbf{Y}_{t+1} | \mathbf{W}_t) = \boldsymbol{\theta}_t + \sum_{l=1}^L \beta_l \mathbf{X}_l + \sum_{r=1}^R m_r(\mathbf{X}_{rt}) + \sum_{s=1}^S m_s(\mathbf{X}_{st}, \mathbf{X}_{size,t}), \quad (1.15)$$

where the notations are similar to [Equation 1.14](#). However, we add a linear term to capture the linear effects of some characteristics, which can increase the rate of convergence and simplify the model and interpretation. Meanwhile, we relax the pairwise interaction between characteristics to a more general form. Similarly, we also assume that both slope parameters and characteristic functions are time-invariant. Therefore, for those nonlinear and interactive characteristics, each characteristic and each pair among them share a certain form of variation.

## 1.6.2 Data Description

Monthly stock returns are collected from CRSP (Center for Research in Security Prices) and security-specific characteristics date is from Compustat. In terms of stock returns, we correct all returns of delisted stock as in [Hou et al. \(2015\)](#). Furthermore, we subtract Fama-French's monthly risk-free returns from monthly stock returns to attain  $\mathbf{Y}$  from July 1967 to June 2017, 600 months in total. As for security-related characteristics matrix  $\mathbf{W}$ , is constructed using the same way of [Freyberger et al. \(2020a\)](#). After trading off the number of assets kept and characteristics' availability, we select 33 characteristics, which are documented in the Appendix. We use balance sheet data ending at fiscal year  $t - 1$  to predict stock excess returns from July  $t - 1$  to June  $t$ . Some characteristics are updated annually, so we take them unchanged during the fiscal year  $t$ . Finally, we merge stock returns and security-specific characteristics.

## 1.6.3 Variable Selection and Model Specification

We apply non-overlapping rolling window analysis in this empirical study. The purpose is to understand whether there are any time variations in [Equation 1.15](#). In each rolling block, we

use pooled panel data to apply S-LASSO. We omit the heterogeneity to assume that the same characteristic has an identical functional form within each rolling window.

For each characteristic, we choose the number of knots to be 6 to construct B-spline bases, which are used to approximate nonlinear effects and choose the number of knots to be 3 for tensor product bases, which are constructed to approximate interactive effects. Next, we substitute all the levels, B-spline bases and tensor products into the S-LASSO algorithm.

There are two steps for S-LASSO to work, and for both steps, similar to simulation studies, we choose  $\lambda_n$  and  $\tilde{\lambda}_n$  through BIC.

We summarize selection results in [Table 1.2](#), [Table 1.3](#) and [Table 1.4](#), respectively. Columns of these tables are rolling window time periods while each row presents selection results of each characteristic separately. We use  $\checkmark$  to show that the corresponding characteristic is selected in a certain rolling block. We omit some rolling blocks due to the non-invertible characteristics matrix. [Table 1.2](#) documents selection results of characteristics' linear effects on assets excess returns. We do not include characteristics that have both linear and nonlinear effects in [Table 1.2](#) as the general effects of these characteristics should be concluded as nonlinear. Compared with [Table 1.3](#), characteristics that only have linear effects on assets returns are uncommon. However, some characteristics experience persistent linear effects on stock returns, such as "C2A" (ratio of cash and short-term investments to total assets), "PCM" (price-to-cost margin), " $r_{12_7}$ " (cumulative past return from 12 to 7 months). [Table 1.2](#) demonstrates that most uni-variate effects from characteristics are nonlinear, and some of them are long-lasting. "LME" (total market capitalization of the previous month), "A2ME" (assets to market capitalization), "AT" (total assets), "E2P" (earnings to price) and "ROA" (return-on-assets) are selected by all rolling windows. Meanwhile, "Investment", "Q" (Tobin's Q), "ROE" (return-on-equity), " $r_{2_1}$ " (short-term reversal 2 to 1 month) and "S2P" (sales-to-price) are frequently chosen. As for interactive effects with firms' sizes, we use "LME" (total market capitalization of the previous month) as the measure of firms' sizes. [Table 1.4](#) shows the characteristics that have interactive effects with "LME". The interactive effects are not limited to be multiplicity by our method. "Free\_cash" is more influential on stock returns when interacting with firms' sizes. "A2ME", "AT", "Q" and "ROA" also substantially interact with "LME".

Empirical results demonstrate the power of S-LASSO to select relevant variables and specify a flexible regression model. We show that asset-related characteristics are relevant to predict stock excess returns. Specifically, the form of each characteristic is different, which includes but is not limited to linear effects, nonlinear effects and interactions with firms' sizes.

Although most uni-variate functions of characteristics are nonlinear, however, linear functions, which have both computational and convergence advantages, are still important. S-LASSO can not only specify linear parts but also select more general interactive effects with firms' sizes since it uses tensor products to approximate more complicated bi-variate functions.

### 1.6.4 Selection Results

Table 1.2 Summary of Linear Effects of Characteristics on Assets Excess Returns

Characteristics	65-68	68-71	71-73	73-76	76-79	79-82	85-88	88-91	91-94	94-97	97-00	03-06	06-09	09-12
LME														
A2ME														
AT														
ATO														
BEME														
C2A		✓	✓		✓						✓		✓	✓
C2D														
CTO	✓		✓	✓	✓									
Delceq														
DelGmSale														
DelshROUT														
E2P														
EPS					✓	✓								
Free_cash	✓		✓	✓										
Investment														
IPM														
Lev	✓				✓	✓								
LTurnover														
PCM	✓	✓	✓		✓	✓								
PM	✓	✓			✓									
Prof														
Q														
ROA														
ROC														
ROE														
r12_2														
r12_7				✓	✓	✓	✓			✓	✓	✓		
r6_2				✓	✓	✓	✓							
r2_1	✓	✓						✓	✓					
S2C														
S2P													✓	
Sales_g			✓											
SGA2S	✓													

This table shows selection results of characteristics that only have linear effects on predicting assets excess returns through three-year rolling windows from July 1965-June 2012. ✓ represents the characteristic is selected in the corresponding rolling window shown in the column.

Table 1.3 Summary of nonlinear Effects of Characteristics on Assets Excess Returns

Characteristics	65-68	68-71	71-73	73-76	76-79	79-82	85-88	88-91	91-94	94-97	97-00	03-06	06-09	09-12
LME	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
A2ME	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
AT	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
ATO	✓				✓	✓		✓	✓	✓	✓	✓	✓	✓
BEME	✓				✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
C2A	✓													
C2D			✓					✓	✓	✓		✓		
CTO							✓	✓	✓	✓	✓	✓	✓	✓
Delceq		✓	✓			✓		✓		✓	✓	✓	✓	✓
DelGmSale	✓		✓					✓	✓	✓		✓	✓	✓
Delshrou			✓		✓		✓		✓	✓	✓	✓	✓	✓
E2P	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
EPS	✓	✓	✓			✓	✓	✓	✓	✓	✓	✓	✓	✓
Free_cash												✓		
Investment	✓				✓		✓	✓	✓	✓		✓	✓	✓
IPM	✓	✓	✓			✓	✓	✓	✓		✓	✓	✓	✓
Lev			✓				✓	✓	✓	✓	✓	✓	✓	✓
LTurnover									✓	✓		✓	✓	✓
PCM										✓		✓	✓	✓
PM					✓		✓	✓		✓	✓	✓	✓	✓
Prof	✓				✓	✓		✓	✓	✓			✓	✓
Q		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
ROA	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
ROC		✓			✓			✓		✓	✓	✓	✓	✓
ROE						✓	✓	✓	✓	✓	✓	✓	✓	✓
r12_2	✓						✓	✓	✓			✓	✓	✓
r12_7								✓				✓		
r6_2							✓	✓	✓	✓	✓	✓	✓	✓
r2_1			✓	✓				✓	✓	✓	✓	✓	✓	
S2C			✓						✓		✓	✓		
S2P			✓		✓		✓	✓	✓	✓	✓	✓	✓	✓
Sales_g			✓	✓				✓		✓		✓		
SGA2S	✓						✓				✓	✓		

This table shows selection results of characteristics that have nonlinear effects on predicting assets excess returns through three-year rolling windows from July 1965-June 2012. ✓ represents the characteristic is selected in the corresponding rolling window shown in the column.

Table 1.4 Summary of Interactive Effects of Characteristics with Size on Assets Excess Returns

Characteristics	65-68	68-71	71-73	73-76	76-79	79-82	85-88	88-91	91-94	94-97	97-00	03-06	06-09	09-12
A2ME	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
AT	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
ATO		✓												
BEME	✓		✓		✓		✓					✓		
C2A	✓		✓		✓		✓		✓	✓	✓	✓		
C2D														
CTO	✓		✓					✓		✓		✓	✓	
Delceq														
DelGmSale	✓													
DelshROUT							✓			✓				
E2P		✓				✓		✓		✓	✓			✓
EPS			✓	✓	✓	✓		✓		✓				
Free_cash					✓			✓	✓	✓		✓	✓	
Investment					✓	✓						✓		
IPM														
Lev					✓	✓			✓	✓	✓	✓		
LTurnover		✓							✓			✓		
PCM			✓		✓	✓			✓	✓		✓		
PM								✓						
Prof						✓								
Q	✓		✓			✓	✓	✓	✓	✓	✓	✓		
ROA	✓	✓	✓		✓				✓	✓	✓	✓	✓	
ROC	✓		✓						✓					
ROE														
r12_2												✓		
r12_7														
r6_2														✓
r2_1									✓	✓				
S2C						✓								
S2P		✓				✓		✓		✓		✓	✓	
Sales_g														
SGA2S			✓					✓	✓			✓		

This table shows selection results of characteristics that have interactive effects with firms' sizes (LME) on predicting assets excess returns through three-year rolling windows from July 1965-June 2012. ✓ represents the characteristic is selected in the corresponding rolling window shown in the column.

## 1.7 Conclusion

We propose a more general variable selection and model specification method, called Specification LASSO (S-LASSO). S-LASSO is designed under sparsity, to specify a partially linear additive non-parametric regression model with pairwise interactions among regressors. Firstly, S-LASSO considers all possibilities through levels, B-splines bases and tensor products of all variables. Then, there are two steps for S-LASSO to work. In the first step, we apply LASSO to give preliminary selection. In the second step, an adaptive group LASSO is employed to give the final selection results in a group manner, using estimates in the first step

---

as discriminatory group penalty. We illustrate the satisfactory accuracy of S-LASSO through simulation studies. Empirically, S-LASSO is applied to a characteristics-based asset pricing model. We show that security-specific characteristics have linear, nonlinear and interactive effects with firms' sizes on assets excess returns, which complements current literature.





## **Chapter 2**

# **Dynamic Peer Groups of Arbitrage Characteristics**



## Abstract

We propose an asset pricing factor model constructed with semiparametric characteristics-based mispricing and factor loading functions. We approximate the unknown functions by B-splines sieve where the number of B-splines coefficients is diverging. We estimate this model and test the existence of the mispricing function by a power enhanced hypothesis test. The enhanced test solves the low power problem caused by diverging B-splines coefficients, with the strengthened power approaches one asymptotically. We also investigate the structure of mispricing components through Hierarchical K-means Clusterings. We apply our methodology to CRSP (Center for Research in Security Prices) and Compustat data for the US stock market with one-year rolling windows during 1967-2017. This empirical study shows the presence of mispricing functions in certain time blocks. We also find that distinct clusters of the same characteristics lead to similar arbitrage returns, forming a "peer group" of arbitrage characteristics.

*Keywords:* Semiparametric; Characteristics-based; Peer Groups; Power-enhanced test

*JEL Classification:* C14; G11; G12

## 2.1 Introduction

Stock returns have both common and firm-specific components. [Ross \(1976\)](#) proposed Arbitrage Pricing Theory (APT) to summarize that expected returns on financial assets can be modeled as a linear combination of risk factors. In such a model, each asset has a sensitivity beta to the risk factor. The APT model explains the excess returns in the cross-sectional direction. [Fama and French \(1993\)](#) and [Fama and French \(2015\)](#) proxied those factors by the returns on portfolios sorted by different characteristics, and they developed three-factor and five-factor models. After extracting the common movement parts, they treated the intercept as the mispricing *alpha*, which is asset-specific and cannot be explained by those risk factors. Many papers use a similar method to present other factor models, such as the four-factor model of [Carhart \(1997\)](#), the q-factor model of [Hou et al. \(2015\)](#), and the factor zoo by [Feng et al. \(2017\)](#) among others. All of the above papers studied observed factors and did not assign characteristics-based information to either alpha or beta.

Security-specific characteristics, such as capitalization and book to market ratio, are usually documented to explain asset-specific excess returns. [Freyberger et al. \(2020a\)](#) analyzed the nonlinear effects of 62 characteristics through Lasso-style regressions. This study concluded that 13 of these characteristics have explanatory power on stock excess returns after selecting by adaptive group Lasso. Characteristics-based information is also exploited to develop arbitrage portfolios by directly parameterizing the portfolio weights as a linear function of characteristics, as in [Hjalmarsson and Manchev \(2012\)](#) and [Kim et al. \(2019\)](#). Empirically, they showed that their portfolio outperformed other baseline competitors.

This paper's contributions are fourfold. Firstly, we build up a more flexible semiparametric characteristics-based asset pricing factor model focusing on the mispricing component. Secondly, we extend previous estimation and testing methods, which can fit the current framework better. Especially, we extend the power-enhanced test of [Fan et al. \(2015\)](#) in a group manner to strengthen the conventional Wald test for mispricing functions. This test can also select the characteristics that contribute to arbitrage portfolios simultaneously. Thirdly, we construct a two-layer clustering structure of mispricing components. Finally, our methods are applied to fifty years of monthly US stock data. We detect distinct clusters of the same characteristics resulting in similar arbitrage returns, forming a "peer group" of arbitrage characteristics. This finding supplements existing portfolio management techniques by implying that the development of arbitrage portfolios through the asset weights determined by the linear mispricing function is improvable.

This class of models has a basic regression specification in [Equation 2.1](#). Consider the panel regression model

$$y_{it} = \alpha_i + \sum_{j=1}^J \beta_{ji} f_{jt} + \varepsilon_{it}, \quad (2.1)$$

where  $y_{it}$  is the excess return of security  $i$  at time  $t$ ;  $f_{jt}$  is the  $j^{\text{th}}$  risk factor's return at time  $t$ ;  $\beta_{ji}$  denotes the  $j^{\text{th}}$  factor loading of asset  $i$ ;  $\alpha_i$  represents the intercept (mispricing) of asset  $i$ ; and  $\varepsilon_{it}$  is the mean zero idiosyncratic shock. In terms of factor loadings  $\beta_{ji}$ , [Connor and Linton \(2007\)](#) and [Connor et al. \(2012\)](#) studied a characteristic-beta model, which bridges the beta-coefficients and firm-specific characteristics by specifying each beta as an unknown function of one characteristic. In their model, beta functions and unobservable factors are estimated by the back-fitting iteration. They concluded that those characteristic-beta functions are significant and nonlinear. Their model can be summarized by

$$y_{it} = \sum_{j=1}^J g_j(X_{ji}) f_{jt} + \varepsilon_{it}, \quad (2.2)$$

where  $X_{ji}$  is the  $j^{\text{th}}$  observable characteristic of firm  $i$ .

They restricted their beta function to be uni-variate and did not consider the components of factor loading functions that cannot be explained by characteristics. To overcome this limitation, [Fan et al. \(2016\)](#) allowed  $\beta_{ji}$  in [Equation 2.2](#) to have a component explained by observable characteristics as well as an unexplained or stochastic part, written as  $\beta_{ji} = g_j(X_i) + u_{ji}$ , where  $u_{ji}$  is mean independent of  $X_{ji}$ . They proposed the Projected Principal Component Analysis (PPCA), which projects stock excess returns onto space spanned by firm-specific characteristics and then applies Principal Component Analysis (PCA) to the projected returns to find the unobservable factors. This method has attractive properties even under the large  $n$  and small  $T$  setting. However, they did not study the mispricing part (alpha), which is crucial to both asset pricing theories and portfolio management.

In this paper, we work on a semiparametric characteristics-based alpha and beta model, which utilizes a set of security-specific characteristics that are similar to [Freyberger et al. \(2020a\)](#). We use unknown multivariate characteristic functions to approximate both  $\alpha_i$  and  $\beta_{ji}$  in [Equation 2.1](#). Specifically, we assume  $\alpha_i$  and  $\beta_{ji}$  are functions of a large set of asset-specific characteristics as  $\alpha_i = h(\mathbf{X}_i) + \gamma_i$  and  $\beta_{ji} = g_j(\mathbf{X}_i) + \lambda_{ij}$ <sup>1</sup>. We only specify additive structure of  $h(\mathbf{X}_i)$  and  $g_j(\mathbf{X}_i)$ , which are further approximated by B-splines sieve. We then estimate

<sup>1</sup> $\mathbf{X}_i$  is a vector of a large set of asset-specific characteristics of stock  $i$ .

$h(\mathbf{X}_i)$ ,  $g_j(\mathbf{X}_i)$  and unobservable risk factors  $f_{jt}$ . In addition, we design a power-enhanced test and Hierarchical K-mean Clustering for the mispricing function  $h(\mathbf{X}_i)$  to study the nonlinear behavior of arbitrage characteristics.

Some recent papers such as [Kim et al. \(2019\)](#) and [Kelly et al. \(2019\)](#) analyzed a similar model as ours, which assume that both  $h(\mathbf{X}_i)$  and  $g_j(\mathbf{X}_i)$  are *linear functions*. They both included around 40 characteristics in  $\mathbf{X}_i$ . However, they drew different conclusions on the existence of  $h(\mathbf{X}_i)$ . [Kim et al. \(2019\)](#) determined assets weights of arbitrage portfolios using one-year rolling window estimates  $\frac{1}{n}\hat{h}(\mathbf{X}_i)$ . They showed that their arbitrage portfolios returns are statistically and economically significant. However, [Kelly et al. \(2019\)](#) applied instrumented principal component analysis (IPCA) to the entire time span from 1965 to 2014, and concluded no evidence to reject the null hypothesis  $H_0 : h(\mathbf{X}_i) = \mathbf{X}_i^\top \mathbf{B} = \mathbf{0}$  through bootstrap. This dispute spurs the development of a more flexible model and reliable hypothesis tests to investigate the existence and structure of  $h(\mathbf{X}_i)$ . The IPCA, which requires both large  $n$  and  $T$  to work, was introduced in [Kelly et al. \(2017\)](#). This method does not fit our setting since we apply rolling window analysis with small  $T$ . Furthermore, [Kelly et al. \(2019\)](#) restricted the function form of  $h(\mathbf{X}_i)$  and  $g_j(\mathbf{X}_i)$  to be time-invariant, which is not consistent with our empirical results under a semiparametric setting. To clarify the differences with the aforementioned research, this paper proposes a semiparametric model, which allows for both nonlinearity and time-variation of  $h(\mathbf{X}_i)$  and  $g_j(\mathbf{X}_i)$ . Furthermore, we consider a different economic question, namely, the existence and structure of mispricing functions. Our empirical study sheds light on why [Kelly et al. \(2019\)](#) and [Kim et al. \(2019\)](#) drew different conclusions: weak, time-varying and nonlinear characteristics-based mispricing functions only appear in certain rolling windows, which is hard to be detected without rolling window analysis. However, given the presence of some persistent arbitrage characteristics, portfolios developed through mispricing functions can provide arbitrage returns.

The unrestrictive model in this paper brings both opportunities and challenges. According to [Huang et al. \(2010b\)](#), the number of B-spline knots must increase in the number of observations to achieve accurate approximation and good asymptotic performance. Therefore, the dimension of B-splines bases coefficients also needs to grow with the sample size. Besides, mispricing functions are treated as anomalies. Under a correctly specified factor model, coefficients of these B-splines bases that are employed to approximate  $h(\mathbf{X}_i)$  are very likely to be sparse. All of these circumstances make the conventional Wald tests have very low power as discussed in [Fan et al. \(2015\)](#). Therefore, a power-enhanced test should be developed to strengthen the power of Wald tests and to detect the most relevant characteristics among a characteristic zoo included in  $h(\mathbf{X}_i)$ . [Kock and Preinerstorfer \(2019\)](#) illustrated that if

the number of coefficients diverges as the number of observations approaches infinity, the standard Wald test is power enhanceable. [Fan et al. \(2015\)](#) proposed a power-enhanced test by introducing a screening process on all estimated coefficients one by one. They added significant components as a supplement to the standard Wald test. In this paper, we extend [Fan et al. \(2015\)](#) to a group manner to enhance the hypothesis test on a high dimensional additive semiparametric function,  $H_0 : h(\mathbf{X}_i) = 0$ . This method allows all the significant components of  $h(\mathbf{X}_i)$  to be selected and contribute to the test statistics, with the test power approaching one.

The careful analysis of  $h(\mathbf{X}_i)$  is theoretically and practically meaningful. Firstly, the presence of  $h(\mathbf{X}_i)$  is an important component of Arbitrage Pricing Theory (APT) and can contribute to asset pricing theories, namely, linking the mispricing functions with security-related characteristics. Secondly, as in [Hjalmarsson and Manchev \(2012\)](#) and [Kim et al. \(2019\)](#),  $h(\mathbf{X}_i)$  can be utilized to construct arbitrage portfolios through observed characteristics. However, both research was built upon the condition that  $h(\mathbf{X}_i)$  is linear over characteristics. If the mispricing function  $h(\mathbf{X}_i)$  is not monotonic, simply setting portfolio weights to the estimated values of linear-specified  $h(\mathbf{X}_i)$  can be problematic. In this paper, we show that some characteristics with substantially different values give rise to similar arbitrage returns. The distance of arbitrage returns between two assets  $i$  and  $j$  is  $d_{ij} = |h(\mathbf{X}_i) - h(\mathbf{X}_j)|$  and the similarity of their characteristics is  $\|\mathbf{X}_i - \mathbf{X}_j\|_2$ , where  $\|\cdot\|_2$  represents  $L_2$  distance. Inspired by [Hoberg and Phillips \(2016\)](#) and [Vogt and Linton \(2017\)](#), we employ a hierarchical K-means clustering to classify these characteristics within each mispricing return group. We present the dynamic of distinct clusters of the same characteristics leading to similar arbitrage returns, forming a "peer group" of arbitrage characteristics. Therefore, under the semiparametric setting, the asset weighting function should rely on these time-varying and nonlinear peer groups.

The rest of this paper is organized as follows. Section 2 sets out the semiparametric model. Section 3 introduces the assumptions and estimation methods. Section 4 constructs a power-enhanced test for high dimensional additive semiparametric functions. Section 5 employs Hierarchical K-Means Clustering to investigate peer groups of arbitrage characteristics. Section 6 describes the asymptotic properties of our estimates and test statistics. Section 6 simulates data to verify the performance of our methodology. Section 7 presents an empirical study. Finally, Section 8 concludes this paper. Characteristics description tables, proofs, mispricing curves and plots of peer groups are arranged in the Appendix.

## 2.2 Model Setup

We assume that there are  $n$  securities observed over  $T$  time periods. We also assume that during a short time window, each security has  $P$  time-invariant observed characteristics, such as market capitalization, momentum, and book-to-market ratios. Meanwhile, we may omit heteroskedasticity by assuming that each characteristic shares a certain form of variation within each period for all securities. We suppose that

$$y_{it} = (h(\mathbf{X}_i) + \gamma_i) + \sum_{j=1}^J (g_j(\mathbf{X}_i) + \lambda_{ij})f_{jt} + \varepsilon_{it}, \quad (2.3)$$

where  $y_{it}$  is the monthly excess return of the  $i^{\text{th}}$  stock at the month  $t$ ;  $\mathbf{X}_i$  is a  $1 \times P$  vector of  $P$  characteristics of stock  $i$  during time periods  $t = 1, \dots, T$ .  $T$  is a small and fixed time block. In practice, most characteristics are updated annually. Thus, we assume  $\mathbf{X}_i$  is time-invariant in one-year time window.  $h(\mathbf{X}_i)$  is an unknown mispricing function explained by a large set of characteristics whereas  $\gamma_i$  is the random intercept of the mispricing part that cannot be explained by characteristics. Similarly, we have characteristics-beta function  $g_j(\cdot)$  to explain the  $j^{\text{th}}$  factor loadings and the unexplained stochastic part of the loading is  $\lambda_{ij}$ .  $\lambda_{ij}$  is orthogonal to the  $g_j(\cdot)$  function.  $f_{jt}$  is the realization of the  $j^{\text{th}}$  risk factor at time  $t$ . Finally,  $\varepsilon_{it}$  is homoskedastic zero-mean idiosyncratic residual of the  $i^{\text{th}}$  stock at time  $t$ . Random variables  $\gamma_i$  and  $\lambda_{ij}$  are used to generalize our settings and not to be estimated. They will be treated as noise in the identification assumptions.

To avoid the curse of dimensionality, we impose additive forms on both  $h(\cdot)$  and  $g_j(\cdot)$  functions:  $h(\mathbf{X}_i) = \sum_{p=1}^P \mu_p(X_{ip})$  and  $g_j(\mathbf{X}_i) = \sum_{p=1}^P \theta_{jp}(X_{ip})$ , where  $\mu_p(X_{ip})$  and  $\theta_{jp}(X_{ip})$  are uni-variate unknown functions of the  $p^{\text{th}}$  characteristic  $X_p$ . We rewrite the model as:

$$y_{it} = \left( \sum_{p=1}^P \mu_p(X_{ip}) + \gamma_i \right) + \sum_{j=1}^J \left( \sum_{p=1}^P \theta_{jp}(X_{ip}) + \lambda_{ij} \right) f_{jt} + \varepsilon_{it}, \quad (2.4)$$

**Assumption 2.2.1.** *We suppose that:*

$$E(\varepsilon_{it} | \mathbf{X}, f_{jt}) = 0,$$

$$E(h(\mathbf{X}_i)) = E(g_j(\mathbf{X}_i)) = 0, \text{ for } j = 1, 2, \dots, J$$

$$E(\gamma_i | \mathbf{X}) = E(\gamma_i),$$

$$E(\lambda_{ij} | \mathbf{X}) = E(\lambda_{ij}), \text{ for } j = 1, 2, \dots, J$$



$$E(h(\mathbf{X}_i)g_j(\mathbf{X}_i)) = \mathbf{0}, \text{ for } j = 1, 2, \dots, J$$

Similar to Connor et al. (2012) and Fan et al. (2016), Assumption 2.2.1 above is to standardize model settings, including the zero mean assumption for factor loadings and mispricing functions for identification purposes. We also impose orthogonality between mispricing and factor loading parts for the identification reason. This is because the variation of risk factors can be absorbed into the mispricing part if it is not orthogonal to the factor loadings. More discussions can be found in Connor et al. (2012).

## 2.3 Estimation

In this section we discuss the approximation of unknown uni-variate functions and our estimation methods for model Equation 2.3. In the semiparametric setting, we apply the Projected-PCA following Fan et al. (2016) to work on the common factors and characteristics-beta directly. Next, we project the residuals onto the characteristics-based alpha space that is orthogonal to the systematic part. The second step is similar to equality-constrained OLS.

### 2.3.1 B-splines Approximation

We use B-splines sieve to approximate unknown functions  $\theta(\cdot)$  and  $\mu(\cdot)$  in Equation 2.4. Similar to Huang et al. (2010b) and Chen and Pouzo (2012), we have the following procedures. Firstly, suppose that the  $p^{th}$  covariate  $X_p$  is in the interval  $[D_0, D]$ , where  $D_0$  and  $D$  are finite numbers with  $D_0 < D$ . Let  $\mathbf{D} = \{\underbrace{D_0, D_0, \dots, D_0}_{l}, d_1 < d_2 < \dots < d_{m_n} < \underbrace{D, D, \dots, D}_{l}\}$  be a simple knot sequence on the interval  $[D_0, D]$ . Here,  $m_n = \lfloor n^\nu \rfloor$  ( $\lfloor \cdot \rfloor$  gives nearest integer) is a positive integer of the number of internal knots, which is a function of security size  $n$  in period  $t$  with  $0 < \nu < 0.5$ .  $l$  is the degree of those bases. Therefore, we have  $H_n = l + m_n$  bases in total, which will diverge as  $n \rightarrow \infty$ . Following this setting, a set of B-splines can be built for the space  $\Omega_n[\mathbf{D}]$ .

Secondly, for the  $p^{th}$  characteristic  $X_p$ , there is a set of  $H_n$  orthogonal bases  $\{\phi_{1p}(X_p), \dots, \phi_{H_n p}(X_p)\}$ . Those uni-variate unknown functions can be approximated as linear combinations of these bases as  $\mu_p(X_p) = \sum_{q=1}^{H_n} \alpha_q \phi_{qp}(X_p) + R_p^\mu(X_p)$  and  $\theta_p(X_p) = \sum_{q=1}^{H_n} \beta_{jq} \phi_{qp}(X_p) + R_p^\theta(X_p)$ , where  $R_p^\mu(X_p)$  and  $R_p^\theta(X_p)$  are approximation errors. It is not necessary to use the same bases for both unknown functions and the representation here is for notational simplicity only. Therefore, the model Equation 2.4 can be written as:

$$y_{it} = \sum_{p=1}^P \left( \sum_{q=1}^{H_n} \alpha_{pq} \phi_{pq}(X_{ip}) + R_p^\mu(X_p) \right) + \gamma_i + \sum_{j=1}^J \left( \sum_{p=1}^P \left( \sum_{q=1}^{H_n} \beta_{jpq} \phi_{pq}(X_{ip}) + R_p^\theta(X_p) \right) + \lambda_{ij} \right) f_{jt} + \varepsilon_{it}$$

For each  $i = 1, 2, \dots, n$ ,  $p = 1, 2, \dots, P$  and  $t = 1, 2, \dots, T$ , we have:

$$\mathbf{1}_T = (1, \dots, 1)^\top \in \mathbb{R}^T,$$

$$\beta_j = (\beta_{1,j1}, \dots, \beta_{H_n,j1}, \dots, \beta_{1,jP}, \dots, \beta_{H_n,jP})^\top \in \mathbb{R}^{H_n P},$$

$$\mathbf{B} = (\beta_1, \dots, \beta_J),$$

$$\mathbf{A} = (\alpha_{11}, \dots, \alpha_{1H_n}, \dots, \alpha_{P1}, \dots, \alpha_{PH_n})^\top \in \mathbb{R}^{H_n P},$$

$$\Phi(\mathbf{X}) = \begin{bmatrix} \phi_{1,11}(X_{11}) & \cdots & \phi_{1,1H_n}(X_{11}) & \cdots & \phi_{1,P1}(X_{1P}) & \cdots & \phi_{1,PH_n}(X_{1P}) \\ \phi_{2,11}(X_{21}) & \cdots & \phi_{2,1H_n}(X_{21}) & \cdots & \phi_{2,P1}(X_{2P}) & \cdots & \phi_{2,PH_n}(X_{2P}) \\ \vdots & \vdots & \vdots & \ddots & \vdots & & \vdots \\ \phi_{n,11}(X_{n1}) & \cdots & \phi_{n,1H_n}(X_{n1}) & \cdots & \phi_{n,P1}(X_{nP}) & \cdots & \phi_{n,PH_n}(X_{nP}) \end{bmatrix},$$

where  $\phi_{i,ph}(X_{ip})$  is the  $h^{\text{th}}$  basis of the  $p^{\text{th}}$  characteristic of asset  $i$  at time  $t$ . Therefore, the original model

$$\mathbf{Y} = (h(\mathbf{X}) + \Gamma) \mathbf{1}_T^\top + (\mathbf{G}(\mathbf{X}) + \Lambda) \mathbf{F}^\top + \mathbf{U},$$

can be represented by B-splines sieve as:

$$\mathbf{Y} = (\Phi(\mathbf{X})\mathbf{A} + \Gamma + \mathbf{R}^\mu(\mathbf{X})) \mathbf{1}_T^\top + (\Phi(\mathbf{X})\mathbf{B} + \Lambda + \mathbf{R}^\theta(\mathbf{X})) \mathbf{F}^\top + \mathbf{U}, \quad (2.5)$$

$\mathbf{Y}$  is the  $n \times T$  matrix of  $y_{it}$ ;  $\Phi(\mathbf{X})$  is the  $n \times PH_n$  matrix of B-splines bases;  $\mathbf{A}$  is a  $PH_n \times 1$  matrix of mispricing coefficients;  $\mathbf{R}^\mu(\mathbf{X})$  is a  $n \times 1$  matrix of approximation errors;  $\mathbf{B}$  is a  $PH_n \times J$  matrix factor loadings' coefficients;  $\mathbf{R}^\theta(\mathbf{X})$  is a  $n \times J$  matrix of approximation errors. We have  $R_p^\mu(X_p) \rightarrow^P 0$  and  $R_p^\theta(X_p) \rightarrow^P 0$ , as  $n \rightarrow \infty$  as in Huang et al. (2010b). Therefore, we omit the approximation errors for simplicity henceforth.  $\mathbf{F}$  is the  $T \times J$  matrix of  $f_{jt}$  and  $\mathbf{U}$  is a  $n \times T$  matrix of  $\varepsilon_{it}$ .  $h(\mathbf{X})$  is a  $n \times 1$  vector of characteristics-based mispricing component;  $\mathbf{G}(\mathbf{X})$  is a  $n \times J$  vector of characteristics-based factor loadings;  $\mathbf{1}_T$  is a  $T \times 1$  vector of 1. The rest are defined the same as Equation 2.4.

We define a projection matrix as:

$$\mathbf{P} = \Phi(\mathbf{X})(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top.$$

The remaining goals of this paper are to estimate both  $h(\mathbf{X})$  and  $\mathbf{G}(\mathbf{X})$  consistently and conduct a power-enhanced test on the hypothesis  $H_0 : h(\mathbf{X}) = \mathbf{0}$ , i.e., to check the existence of mispricing functions under semiparametric settings. Finally, we cluster peer groups of arbitrage characteristics.

### 2.3.2 Two-Step Projected-PCA

In this section, we combine and extend Projected-PCA by Fan et al. (2016) and equality constrained least squares similar to Kim et al. (2019) to estimate the model. To facilitate the estimation, we define a  $T \times T$  time series demeaning matrix  $\mathbf{D}_T = \mathbf{I}_T - \frac{1}{T} \mathbf{1}_T \mathbf{1}_T^\top$ .<sup>2</sup> Next, we demean the equation above on both sides. Therefore we have

$$\mathbf{YD}_T = \tilde{\mathbf{Y}} = (\Phi(\mathbf{X})\mathbf{B} + \Lambda)\mathbf{F}^\top \mathbf{D}_T + \mathbf{UD}_T.$$

Mispricing terms disappear since they are time-invariant by  $(\Phi(\mathbf{X})\mathbf{A} + \Gamma)\mathbf{1}_T^\top \mathbf{D}_T = \mathbf{0}$ . This helps us to work on the systematic part later. Henceforth, we use  $\mathbf{F}$  to represent the time-demeaned factor matrix.

Our procedures are designed to estimate factor loadings  $\mathbf{G}(\mathbf{X})$ , time-demeaned unobserved factors  $\mathbf{F}$  and mispricing coefficients  $\mathbf{A}$  in sequence.

Under Assumption 2.2.1, we have the following estimation procedures:

- 1 Projecting  $\tilde{\mathbf{Y}}$  onto the spline space spanned by  $\{\mathbf{X}_{ip}\}_{i \leq n, p \leq P}$  through a  $n \times n$  projection matrix  $\mathbf{P}$  with  $\mathbf{P} = \Phi(\mathbf{X})(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top$ . We then collect the projected data  $\hat{\mathbf{Y}} = \Phi(\mathbf{X})(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top \tilde{\mathbf{Y}}$ .
- 2 Applying the Principal Component Analysis to the projected data  $\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$ . This allows us to work directly on the sample covariance of  $\mathbf{G}(\mathbf{X})\mathbf{F}^\top$ , under the condition  $E(g_j(\mathbf{X}_i)\varepsilon_{it}) = E(g_j(\mathbf{X}_i)\lambda_{ij}) = 0$ .
- 3 Estimating  $\hat{\mathbf{F}}$  as the eigenvectors corresponding to the first  $J$  (assumed given) eigenvalues of the  $T \times T$  matrix  $\frac{1}{n} \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$  (covariance of projected  $\hat{\mathbf{Y}}$ ).

The method above substantially improves estimation accuracy and facilitates theoretical analysis even under the large  $n$  and small  $T$ . Small  $T$  is preferable in our model setting as we use one-year rolling windows analysis in both simulation and empirical studies, and large  $n$  is required for asymptotic analysis.

---

<sup>2</sup> $\mathbf{I}_T$  is a  $T \times T$  identity matrix, and  $\mathbf{1}_T$  is a  $T \times 1$  matrix of 1.

Factor loadings  $\hat{\mathbf{G}}(\mathbf{X})$  are estimated as:

$$\hat{\mathbf{G}}(\mathbf{X}) = \hat{\mathbf{Y}}\hat{\mathbf{F}}(\hat{\mathbf{F}}^\top\hat{\mathbf{F}})^{-1}$$

In the next step, we estimate the coefficients of the mispricing bases.

4 The estimator of  $\mathbf{A}$  is

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \text{vec}(\mathbf{Y} - \Phi(\mathbf{X})\mathbf{A}\mathbf{1}_T^\top - \hat{\mathbf{G}}(\mathbf{X})\hat{\mathbf{F}}^\top)^\top \text{vec}(\mathbf{Y} - \Phi(\mathbf{X})\mathbf{A}\mathbf{1}_T^\top - \hat{\mathbf{G}}(\mathbf{X})\hat{\mathbf{F}}^\top), \quad (2.6)$$

subject to  $\hat{\mathbf{G}}(\mathbf{X})^\top\Phi(\mathbf{X})\mathbf{A} = \mathbf{0}_J$ .

Let a  $PH_n \times 1$  vector  $\hat{\mathbf{A}}$  be a closed-form solution:

$$\hat{\mathbf{A}} = \mathbf{Q}\tilde{\mathbf{A}},$$

where

$$\mathbf{Q} = \mathbf{I} - (\Phi(\mathbf{X})^\top\Phi(\mathbf{X}))^{-1}\Phi(\mathbf{X})^\top\hat{\mathbf{G}}(\mathbf{X})(\hat{\mathbf{G}}(\mathbf{X})^\top\hat{\mathbf{G}}(\mathbf{X}))^{-1}\hat{\mathbf{G}}(\mathbf{X})^\top\Phi(\mathbf{X}),$$

$$\tilde{\mathbf{A}} = \frac{1}{T}(\Phi(\mathbf{X})^\top\Phi(\mathbf{X}))^{-1}\Phi(\mathbf{X})^\top(\mathbf{Y} - \hat{\mathbf{G}}(\mathbf{X})\hat{\mathbf{F}}^\top)\mathbf{1}_T,$$

given  $\mathbf{P}\hat{\mathbf{G}}(\mathbf{X}) = \hat{\mathbf{G}}(\mathbf{X})$ .

As in Assumption 2.2.1, the  $h(\mathbf{X})$  is orthogonal to the characteristics-based loadings  $\mathbf{G}(\mathbf{X})$ .

5 We also estimate the covariance matrix of  $\hat{\mathbf{A}}$ , i.e.,  $\hat{\Sigma}$ , by extending the methods of Liew (1976). This can facilitate theoretical analysis in the next section. According to Liew (1976),  $\hat{\mathbf{A}}$  is the equality constrained least-square estimator, which has the covariance matrix as (under  $n \leq T$  and covariance shrinkage as in Ledoit et al. (2012) and Fan et al. (2013) among others.):

$$\hat{\Sigma} = \mathbf{Q}\hat{\Sigma}_{\hat{\mathbf{A}}}\mathbf{Q}^\top,$$

where:

$$\hat{\Sigma}_{\hat{\mathbf{A}}} = (\Phi(\mathbf{X})^\top\Phi(\mathbf{X}))^{-1}\Phi(\mathbf{X})^\top \begin{bmatrix} \hat{\sigma}_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{\sigma}_n^2 \end{bmatrix} \Phi(\mathbf{X})(\Phi(\mathbf{X})^\top\Phi(\mathbf{X}))^{-1},$$

$$\hat{\sigma}_i^2 = \frac{\sum_1^T \hat{e}_{it}^2}{T-1},$$

where  $\sum_1^T \hat{e}_{it}^2 = \sum_1^T (y_{it} - \bar{y}_i - \sum_{p=1}^P \sum_{q=1}^{H_n} \hat{\alpha}_{pq} \phi_{pq}(x_{ip}) - \sum_{j=1}^J (\sum_{p=1}^P \sum_{q=1}^{H_n} \hat{\beta}_{jpq} \phi_{pq}(x_{ip})) \hat{f}_{jt})^2$ .  
Heteroskedasticity is caused by  $\gamma_i$ .

## 2.4 Power-enhanced Tests

There are considerable discussions about the mispricing phenomenon under factor models, while the existence of mispricing functions remains controversial. Namely, whether there are relevant covariates explaining remaining excess returns after subtracting co-movements components captured by risk factors. Recently, [Kim et al. \(2019\)](#) found the characteristics arbitrage opportunities by estimating a linear characteristic mispricing function without providing theoretical results. However, [Kelly et al. \(2019\)](#) conducted a conventional Wald hypothesis test on the similar mispricing function using bootstrap, concluding that there is no evidence to reject the null hypothesis  $H_0 : h(\mathbf{X}) = \mathbf{0}$ . Additionally, they applied the bootstrap method to estimate the covariance matrix  $\Sigma$ , which caused potential problems for theoretical analysis. Moreover, according to [Fan et al. \(2015\)](#), their test results may have relatively low power when the true coefficient vector of linear mispricing function  $\mathbf{A}$  has a sparse structure.

Both studies adopt a parametric framework, which relies on the strong assumption of linearity. However, this assumption is not consistent with [Connor et al. \(2012\)](#), which showed that both characteristic-beta and mispricing functions are very likely to be nonlinear. Therefore, we propose a semiparametric model to accommodate the nonlinearity to a great extent.

But semiparametric framework leads to additional challenges for inference. On the one hand, as mentioned above, the number of coefficients of mispricing B-splines diverges as  $n \rightarrow \infty$ , which implies that the power of the standard Wald test can be quite low, (see [Fan et al. \(2015\)](#)). On the other hand, according to other research like [Fama and French \(1993\)](#) and [Fama and French \(2015\)](#), mispricing terms can be regarded as anomalies. This means that in our model setting, the true mispricing coefficient vector  $\mathbf{A}$  can be high-dimensional but sparse, reducing the power of the conventional Wald test further.

As in [Kock and Preinerstorfer \(2019\)](#), conventional hypothesis tests under these circumstances are power enhanceable. The power-enhanced Wald test in this paper is an extension of [Fan et al. \(2015\)](#) to a group manner, namely, the hypothesis test under high-dimensional additive semiparametric settings. The proposed test is power strengthened when the dimension of the coefficients of the additive regression  $\mathbf{A}$  is diverging as  $n \rightarrow \infty$  without size distortion.

Meanwhile, this test is robust to sparse alternatives. On top of that, the proposed test can select the most important components from sparse additive functions. Finally, the proposed method can also be applied when the number of characteristics is diverging, i.e.,  $P \rightarrow \infty$ .

We construct a new test:

$$H_0 : h(\mathbf{X}) = \mathbf{0}, \quad H_1 : h(\mathbf{X}) \neq \mathbf{0},$$

equivalently,

$$H_0 : \mathbf{A} = \mathbf{0}, \quad H_1 : \mathbf{A} \in \mathcal{A},$$

where  $\mathcal{A} \subset \mathbb{R}^{PH_n} \setminus \mathbf{0}$ .

Here, we have:

$$S_1 = \frac{\hat{\mathbf{A}}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{A}} - PH_n}{\sqrt{2PH_n}},$$

where  $S_1$  is the "original" Wald test statistics;  $P$  is the number of characteristics;  $PH_n$  is the total number of B-spline bases, and  $\hat{\mathbf{A}} \in \mathbb{R}^{PH_n}$ . The value of  $H_n$  is a function of asset number  $n$ , therefore,  $H_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Under  $H_0$ ,  $S_1$  has a nondegenerate limiting distribution  $F$  as  $n \rightarrow \infty$ . Given the significance level  $q$ ,  $q \in (0, 1)$  as well as the critical value  $F_q$ :

$$S_1 | H_0 \rightarrow^d F,$$

$$\lim_{n \rightarrow \infty} \Pr(S_1 > F_q | H_0) = q.$$

[Pesaran and Yamagata \(2012\)](#) showed that:

$$S_1 | H_0 \rightarrow^d \mathcal{N}(0, 1),$$

under regularity conditions.

Potentially, sparse and diverging  $PH_n$  means that it is plausible to add a power-enhanced component to  $S_1$ , which can improve the power of the hypothesis test without any size distortions.

Therefore, we can construct an extra screening component  $S_0$  as:

$$S_0 = H_n \sum_{p=1}^P \mathbf{I} \left( \sum_{h=1}^{H_n} |\hat{\alpha}_{ph}| / \hat{\sigma}_{ph} \geq \eta_n \right),$$

where  $\hat{\sigma}_{ph}$  is the square-root of the  $ph^{th}$  entry of the diagonal elements of  $\hat{\Sigma}$ .  $\mathbf{I}(\cdot)$  is an indicator for the screening process while  $\eta_n$  is a data-driven threshold value to avoid potential size-distortion.

Here we discuss the choice of  $\eta_n$ . By construction and the assumption of independent characteristics, we assume that all B-splines bases are orthogonal. Our goal is to bound the maximum of those standardized coefficients under the null hypothesis.

Define  $Z = \max_{\{1 \leq p \leq P, 1 \leq h \leq H_n\}} \{|\hat{\alpha}_{ph}|/\hat{\sigma}_{ph}\}$ . We have

$$\hat{\alpha}_{ph}/\hat{\sigma}_{ph}|\mathbf{H}_0 \rightarrow^d N(0, 1),$$

$$E(Z|\mathbf{H}_0) = \sqrt{2 \log PH_n}.$$

After grouping the coefficients of bases used to approximate the unknown function of each characteristic, let  $R = \max(\sum_{h=1}^{H_n} |\hat{\alpha}_{1h}|/\hat{\sigma}_{1h}, \dots, \sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph}, \dots, \sum_{h=1}^{H_n} |\hat{\alpha}_{Ph}|/\hat{\sigma}_{Ph})$ . Following this, we set the threshold as  $\eta_n = H_n \sqrt{2 \log(PH_n)}$ , where  $H_n = l + n^\nu$ . As  $H_n$  is a diverging sequence, it can control the influence of the group size properly. Meanwhile,  $\eta_n$  also diverges so that  $\eta_n$  is a conservative threshold value used to avoid potential size distortion.

Apart from strengthening the power of conventional hypothesis test,  $\mathbf{I}(\cdot)$  is a screening term that can select the most relevant characteristics at the same time.

We then define the arbitrage characteristics set, which includes the characteristics that have the strong explanation power for mispricing functions:

$$\hat{\mathcal{M}} = \{\mathbf{X}_p \in \mathbf{X} : \sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n, \quad p = 1, 2, \dots, P\},$$

and  $M$  is the cardinality of the set containing mispricing characteristics. When the set  $\hat{\mathcal{M}}$  is relatively small, conventional tests are likely to suffer the lower power problem. The added  $S_0$  strengthens the power of the test and drives the power to one since  $S_0$  is slowly diverging.

Therefore, our new test statistic is  $S = S_0 + S_1$ , and asymptotic properties of  $S$  will be discussed later.

To conclude, the advantages of  $S = S_0 + S_1$  are:

- 1 The power of the hypothesis test on  $H_0 : h(\mathbf{X}) = \mathbf{0}$  is mainly enhanced without size distortions.
- 2 We can find specific characteristics which cause the mispricing by this screening mechanism.

As designed,  $S_0$  satisfies all three properties of [Fan et al. \(2015\)](#), as  $n \rightarrow \infty$ :

- 1  $S_0$  is non-negative,  $\Pr(S_0 \geq 0) = 1$
- 2  $S_0$  does not cause size distortion: under  $H_0$ ,  $\Pr(S_0 = 0 | H_0) \rightarrow 1$
- 3  $S_0$  enhances test power. Under  $H_1$ ,  $S_0$  diverges quickly in probability given the well chosen  $\eta_n$ .

Based on properties of  $S_0$ , we have three properties of  $S$  listed:

- 1 No size-distortion  $\limsup_{n \rightarrow \infty} \Pr(S > F_q | H_0) = q$
- 2  $\Pr(S > F_q | H_1) \geq \Pr(S_1 > F_q | H_1)$ . Hence, the power of  $S$  is at least as large as that of  $S_1$ .
- 3  $\Pr(S > F_q | \hat{\mathcal{M}} \neq \emptyset) \rightarrow 1$  when  $S_0$  diverges. This happens, especially, when the true form of  $\mathbf{A}$  has a sparse structure.

## 2.5 Hierarchical K-Means Clustering

This section introduces a Hierarchical K-means Clustering method to find peer groups of arbitrage characteristics based on their arbitrage returns. We ask whether distinct groups of the same arbitrage characteristics, according to their similarity measured by  $\|\mathbf{X}_i - \mathbf{X}_j\|_2$ , may result in similar characteristic-based arbitrage returns in each rolling block, which is an implication for the non-monotonic mispricing function, and forms a "peer group" of arbitrage characteristics. Because traditional arbitrage portfolios as in [Kim et al. \(2019\)](#) and [Hjalmarsson and Manchev \(2012\)](#) rely on the linearity of characteristics-based mispricing components to work, our clustering results can show whether this method is still applicable under more flexible semiparametric model. If there are persistent peer groups in arbitrage returns, investors should consider to long the assets in the peer group with the highest



arbitrage returns while short the assets in the peer group with the lowest arbitrage returns to form an arbitrage portfolio.

Introduction of K-means clustering can be found in [Cox \(1957\)](#) and [Fisher \(1958\)](#).

After the screening process in [section 2.4](#), we obtain the relevant components of mispricing function  $h(\mathbf{X})$ , which is estimated as

$$\hat{\mathcal{M}} = \{\mathbf{X}_p \in \mathbf{X} : \sum_{h=1}^{H_n} |\hat{\alpha}_{ph}| / \hat{\sigma}_{ph} \geq \eta_n, \quad p = 1, 2, \dots, P\}.$$

We define an  $n \times M$  matrix  $\mathbf{M}$  of arbitrage characteristics at time window  $t$  as :

$$\mathbf{M} = \{\mathbf{X}_1, \dots, \mathbf{X}_m, \dots, \mathbf{X}_M\}, \text{ where } \mathbf{X}_m \in \hat{\mathcal{M}}.$$

Note that these characteristics are time-invariant in each rolling window. We also set characteristics-based arbitrage returns of asset  $i$  in month  $t$  as:

$$\ddot{y}_{it} = \phi(\mathbf{M}_i) \hat{\mathbf{A}}_M,$$

where  $\phi(\mathbf{M}_i)$  and  $\hat{\mathbf{A}}_M$  are the corresponding parts of matrix  $\Phi(\mathbf{X}_i)$  and vector  $\hat{\mathbf{A}}$ . For each rolling window, we classify all  $n$  assets through a 2-layer K-means clustering. At the first layer, we cluster these assets into  $K$  groups according to the similarity of their characteristics-based arbitrage returns  $\ddot{y}_{it}$ . At the second layer, we divide  $R_j$  subgroups within the  $j^{\text{th}}$  group from the first layer by the similarity of their arbitrage characteristics, where  $j = 1, 2, \dots, K$ . Finally, the peer groups of arbitrage characteristics can be attained. We repeat this method for all rolling blocks to investigate dynamic patterns of these peer groups. These clusterings will provide illustrative evidence of linear/nonlinear and time-invariant/time-varying structure of mispricing function  $h(\mathbf{X})$ .

We give the classification procedures of both layers. We define  $\Delta_{ij}$  as the difference between characteristics-based arbitrage returns of  $\ddot{y}_{it}$  and  $\ddot{y}_{jt}$ , as well as  $\Upsilon_{ij}$  as the difference between arbitrage characteristics:

$$\Delta_{ij} = \ddot{y}_{it} - \ddot{y}_{jt}, \text{ where } i \neq j, \quad i, j = 1, 2, \dots, n.$$

$$\Upsilon_{ij} = \|\mathbf{M}_i - \mathbf{M}_j\|_2, \text{ where } i \neq j, \quad i, j = 1, 2, \dots, n,$$

$\mathbf{M}_i$  represents the  $i^{th}$  row of  $\mathbf{M}$ . We set two tolerance thresholds  $\psi_y$  and  $\psi_x$ , which are used to control the biggest difference within each group of both layers separately. To accelerate the convergence of the K-means Clustering, we first apply a first difference process, which is introduced below, to obtain centroids as in [Vogt and Linton \(2017\)](#).

For the first layer, we have first difference process:

1. **First difference:** We randomly pick  $i^{th}$  asset and then we calculate  $\Delta_{ij}$  with other assets  $j = 1, 2, \dots, n$ . Thus we obtain  $\Delta_{i(1)} \dots \Delta_{i(n)}$ , with  $n$  being the total individuals for classification. Without loss of generality, we assume  $\Delta_{i(1)} = \min\{\Delta_{i(1)}, \dots, \Delta_{i(n)}\}$ , and  $\Delta_{i(n)} = \max\{\Delta_{i(1)}, \dots, \Delta_{i(n)}\}$ .
2. **Ordering:** We rank the values obtained in Step 1 as follows:

$$\begin{aligned} \Delta_{i(1)} &\leq \dots \leq \Delta_{i(j_1-1)} < \Delta_{i(j_1)} \leq \dots \leq \Delta_{i(j_2-1)} \\ &< \Delta_{i(j_2)} \leq \dots \leq \Delta_{i(j_3-1)} \\ &\vdots \\ &< \Delta_{i(j_{K-1})} \leq \dots \leq \Delta_{i(n)}. \end{aligned}$$

We use the strict inequality mark to show large jumps of "first difference", all of which are larger than  $\psi_y$ , while the weak inequality means that the distance calculated is not larger than  $\psi_y$ . We identify  $K - 1$  jumps that are larger than  $\psi_y$  above. Thus, the initial classification is achieved, and we have a total of  $K$  groups with  $j_1 - 1$  members in the first group  $\mathcal{C}_1$ ,  $j_2 - j_1$  members in the second group  $\mathcal{C}_2$ , ..., and  $n - j_{K-1} + 1$  members in the final group  $\mathcal{C}_K$ .

In terms of the second layer, for the assets in the  $k^{th}$  group  $\mathcal{C}_k$ , we use the same method to further divide them into  $r$  subgroups as  $\mathcal{R}_{1k}, \mathcal{R}_{2k}, \dots, \mathcal{R}_{rk}$ . Within each subgroup, we have:

$$\Upsilon_{ab} = \|\mathbf{M}_a - \mathbf{M}_b\|_2 \leq \psi_x, \text{ where } a, b \in \mathcal{R}_{ik}, i = 1, 2, \dots, r, \text{ and } k = 1, 2, \dots, K.$$

The K-means algorithm is:

1. Step 1: Determine the starting mean values for each group  $\hat{c}_1^{[0]}, \dots, \hat{c}_K^{[0]}$  and calculate the distances  $\hat{D}_k(i) = \Delta(\dot{y}_{it}, \hat{c}_k^{[0]}) = |\dot{y}_{it} - \hat{c}_k^{[0]}|$  for each  $i$  and  $k$ . Define the partition  $\{\mathcal{C}_1^{[0]}, \dots, \mathcal{C}_K^{[0]}\}$  by assigning the  $i^{th}$  individual to the  $k$ -th group  $\mathcal{C}_k^{[0]}$  if  $\hat{D}_k(i) = \min_{1 \leq k' \leq K} \hat{D}_{k'}(i)$ .

2. Step  $l$ : Let  $\{\mathcal{C}_1^{[l-1]}, \dots, \mathcal{C}_K^{[l-1]}\}$  be the partition of  $\{1, \dots, n\}$  from the latest iteration step. Calculate mean functions

$$\hat{c}_k^{[l]} = \frac{1}{|\mathcal{C}_k^{[l-1]}|} \sum_{i \in \mathcal{C}_k^{[l-1]}} \ddot{y}_{it} \quad \text{for } 1 \leq k \leq K$$

And then we calculate  $\Delta(\ddot{y}_{it}, \hat{c}_k^{[l]}) = |\ddot{y}_{it} - \hat{c}_k^{[l]}|$  for each  $i$  and  $k$ . Define the partition  $\{\mathcal{C}_1^{[l]}, \dots, \mathcal{C}_K^{[l]}\}$  by assigning the  $i^{\text{th}}$  individual to the  $k$ -th group  $\mathcal{C}_k^{[l]}$  if  $\hat{D}_k(i) = \min_{1 \leq k' \leq K} \hat{D}_{k'}(i)$ .

3. Iterate the above steps until the partition  $\{\mathcal{C}_1^{[w]}, \dots, \mathcal{C}_K^{[w]}\}$  does not change anymore.

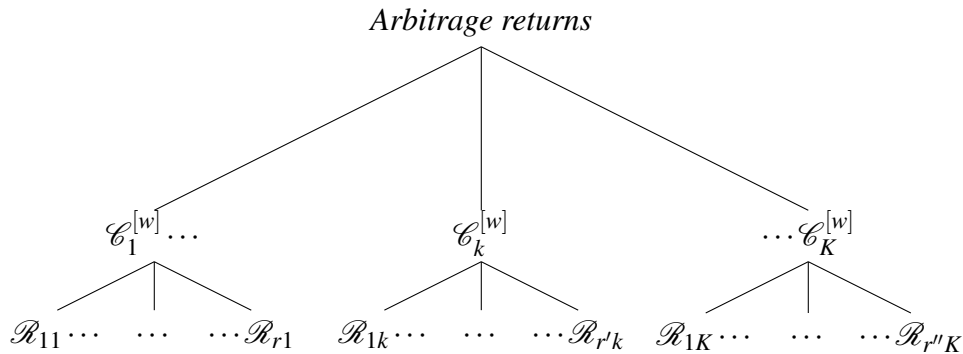
To accelerate the convergence of K-means algorithm, at the step 1, results of first difference are used. As we have already obtained our initial grouping  $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ , therefore starting values for the Step 1 is:

$$\hat{c}_k^{[0]} = \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} \ddot{y}_{it} \quad \text{for } 1 \leq k \leq K,$$

where  $|\mathcal{C}_k|$  is the cardinality of the group  $\mathcal{C}_k$ .

The consistency and other theoretical results of the above procedures can be found in [Pollard \(1981\)](#), [Pollard et al. \(1982\)](#), [Sun et al. \(2012\)](#) and [Vogt and Linton \(2017\)](#).

For the second layer, we repeat the procedures within each group  $\mathcal{C}_k^{[w]}$  with respect to  $\Upsilon_{ab}$ , and the structure of characteristics-based arbitrage returns is:



The first layer is the structure of characteristics-based arbitrage returns, while the second layer gives peer groups of characteristics that can provide similar characteristics-based arbitrage returns.

The number of clusterings is determined by threshold values  $\psi_y$  and  $\psi_x$  directly.  $\psi_y$  and  $\psi_x$  are chosen by the trade-off between the number of clusterings and total within-group sum of squares.

## 2.6 Asymptotic Properties

This section discusses assumptions and properties of estimates and power enhanced statistics  $S$ .

**Definition 2.6.1.** We define  $\mathbf{A} \rightarrow^P \mathbf{B}$  as  $n \rightarrow \infty$  of two  $n \times m$  matrix  $\mathbf{A}$  and  $\mathbf{B}$  with fixed  $p$  when  $\frac{1}{n}(\mathbf{A} - \mathbf{B})^\top(\mathbf{A} - \mathbf{B}) \rightarrow^P \mathbf{0}_{m \times m}$  as  $n \rightarrow \infty$ .

### 2.6.1 Consistency Assumptions

**Assumption 2.6.1.** As  $n \rightarrow \infty$ , we have:

$$\frac{1}{n} \mathbf{Y}^\top \mathbf{Y} \rightarrow^P \mathbf{M}_Y,$$

$$\mathbf{F}^\top \mathbf{F} = \mathbf{I}_J,$$

where  $\mathbf{M}_Y$  is a positive definite matrix, and  $\mathbf{I}_J$  is a  $J \times J$  identity matrix.

We define  $\lambda_{\min}(M)$  and  $\lambda_{\max}(M)$  as the largest and the smallest eigenvalues of matrix  $M$ , respectively. Additionally, we define  $C_{\min}$  and  $C_{\max}$  to be positive constants such that:

$$C_{\min} \leq \lambda_{\min}\left(\frac{1}{n} \Phi^\top(\mathbf{X}) \Phi(\mathbf{X})\right) < \lambda_{\max}\left(\frac{1}{n} \Phi^\top(\mathbf{X}) \Phi(\mathbf{X})\right) \leq C_{\max}$$

as  $n \rightarrow \infty$ .

We impose these restrictions to avoid non-invertibility of stock returns, characteristics, and rotation indeterminacy.

**Assumption 2.6.2.**

$$\frac{1}{n} \mathbf{G}(\mathbf{X})^\top \mathbf{G}(\mathbf{X}) \rightarrow^P \begin{bmatrix} d_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & d_J \end{bmatrix},$$

as  $n \rightarrow \infty$ , where  $d_j$  are distinct and positive entries.

Both Assumption 2.6.1 and 2.6.2 are similar to those in Fan et al. (2016), which are used to separately identify risk factors and factor loadings. Given the orthogonal bases of B-splines and uncorrelated or weakly correlated characteristics, Assumption 2.6.2 is mild.

**Assumption 2.6.3.**  $K_{min}$  and  $K_{max}$  are positive constants such that:

$$K_{min} \leq \lambda_{min}\left(\frac{1}{n}\mathbf{G}(\mathbf{X})^\top \mathbf{P}\mathbf{G}(\mathbf{X})\right) < \lambda_{max}\left(\frac{1}{n}\mathbf{G}(\mathbf{X})^\top \mathbf{P}\mathbf{G}(\mathbf{X})\right) \leq K_{max}$$

as  $n \rightarrow \infty$ .

This assumption requires non-vanishing explanatory power of the B-splines bases  $\Phi(\mathbf{X})$  on the factor loading matrix  $\mathbf{G}(\mathbf{X})$ .

**Assumption 2.6.4.**  $\varepsilon_{it}$  is realized i.i.d. idiosyncratic shocks with  $E(\varepsilon_{it}) = 0$  and  $\text{var}(\varepsilon_{it}) = \sigma^2$ .

Heteroskedasticity is caused by  $\gamma_i$ , namely,  $\text{var}(\gamma_i + \varepsilon_{it}) = \sigma_i^2$ .

## 2.6.2 Main Results

**Theorem 2.6.1.** Let  $\hat{\mathbf{F}}$  be the  $T \times J$  matrix estimate of latent risk factors. Under Assumption 2.2.1-2.6.3,  $\hat{\mathbf{F}} \rightarrow^P \mathbf{F}$ , as  $n \rightarrow \infty$ .

**Theorem 2.6.2.** Define the  $n \times J$  matrix  $\hat{\mathbf{G}}(\mathbf{X})$  as the estimate of factor loadings  $\mathbf{G}(\mathbf{X})$ . Under Assumption 2.2.1-2.6.3 and Theorem 2.6.2, as  $n \rightarrow \infty$ , then  $\hat{\mathbf{G}}(\mathbf{X}) \rightarrow^P \mathbf{G}(\mathbf{X})$ .

**Theorem 2.6.3.** Let the  $PH_n \times 1$  vector  $\hat{\mathbf{A}}$  be the solution of constrained OLS Equation 2.6, then

$$\hat{\mathbf{A}} = \mathbf{Q}\tilde{\mathbf{A}},$$

where

$$\mathbf{Q} = \mathbf{I} - (\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}) (\hat{\mathbf{G}}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}))^{-1} \hat{\mathbf{G}}(\mathbf{X})^\top \Phi(\mathbf{X}),$$

$$\tilde{\mathbf{A}} = \frac{1}{T} (\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top (\mathbf{Y} - \hat{\mathbf{G}}(\mathbf{X})\hat{\mathbf{F}}^\top) \mathbf{1}_T^\top.$$

Under Assumption 2.2.1-2.6.3,  $\Phi(\mathbf{X})\hat{\mathbf{A}} \rightarrow^P h(\mathbf{X})$ , as  $n \rightarrow \infty$ .

**Theorem 2.6.4.** Under Assumption 2.6.2 and Assumption 2.6.4,  $E(Z|\mathbf{H}_0) = \sqrt{2\log PH_n}$ .

**Theorem 2.6.5.** Under  $n \rightarrow \infty$  and  $H_0$ , given the properties of  $S_0$  and  $S_1$ , then:

$$S \rightarrow^d N(0,1),$$

the power of  $S$  is approaching 1 once the arbitrage characteristic is selected as:

$$\Pr(\text{reject } H_0 | \hat{\mathcal{M}} \neq \emptyset) \rightarrow 1.$$

## 2.7 Numerical Study

In this section, we use Compustat and Fama-French three and five factors data to simulate stock returns and then evaluate the performance of our estimation and hypothesis test procedures.

### 2.7.1 Data Generation

Firstly, we use Fama-French three factors monthly returns and all the characteristics that will be included in the empirical study to simulate the stock excess returns. Most of the characteristics are updated annually so we treat those variables as time-invariant during each one-year rolling block. For the characteristics that are updated every month, we substitute the mean values as their fixed values for each fiscal year. We use Fama-French monthly returns from July of year  $t$  to June of year  $t + 1$  and characteristics of fiscal year  $t - 1$  to generate the stock returns from July of year  $t$  to June of year  $t + 1$ . The periods we generate are the same as the empirical study, namely, 50 years from July 1967 to June 2017. For each rolling block with 12 months we have:

$$y_{it} = h(X_i) + \sum_{j=1}^3 g_j(\mathbf{X}_{ij})f_{jt} + \varepsilon_{it}, \quad (2.7)$$

$y_{it}$  is the generated stock's return;  $h(X_i)$  is the mispricing function consisting of a nonlinear characteristic function of  $x_i$ , which is to mimic the sparse structure of the mispricing function;  $g_j(\mathbf{X}_{ij})$  is the  $j^{\text{th}}$  characteristics-based factor loading, which has an additive semiparametric structure;  $\mathbf{X}_{ij}$  is the  $j^{\text{th}}$  subset consisting of 4 characteristics;  $f_{jt}$  is the  $j^{\text{th}}$  Fama-French factor returns at time  $t$ ;  $\varepsilon_{it}$  is the idiosyncratic shock for stock  $i$  at time  $t$ , generated from  $N(0, \sigma^2)$ .

We generate characteristic functions:

$$h(X_i) = \sin X_i,$$

$$g_j(\mathbf{X}_{ij}) = X_{ij1}^2 + (3X_{ij2}^3 - 2X_{ij2}^2) + (3X_{ij3}^3 - 2X_{ij3}^2) + X_{ij4}^2,$$

$X_{ijl}$  is a randomly picked characteristic without replacement from the data in empirical study and  $j = 1, 2, 3, l = 1, \dots, 4$ . A description of these characteristics can be found in the Appendix. Additionally, all  $h(X_i)$ ,  $g_j(\mathbf{X}_{ij})$  are rescaled to have zero mean and unit variance. As we use real data to conduct the simulation, the Assumption of independent  $X_i$  may not be satisfied. Although some characteristics are correlated, the semiparametric model overcomes this problem properly when compared with the parametric model that has serious size distortion.

We do not specify  $h(X_i)$  and  $g_j(\mathbf{X}_{ij})$  to be orthogonal explicitly, but we draw characteristics without replacement and employ sine-waves and polynomials to approximate the orthogonality as much as possible. In this simulation, our method can only estimate the component of  $h(X_i)$  that is orthogonal to  $g_j(\mathbf{X}_{ij})$ . However, results reveal that one can still select the arbitrage characteristics even if we cannot estimate arbitrary  $h(X_i)$  unrestrictedly.

### 2.7.2 Model Misspecification

In this subsection, we show the necessity to consider semiparametric analysis when the forms of factor loadings and mispricing functions are nonlinear.

Under the data generation process, we consider both semiparametric and linear analysis to compare Mean Squared Error (MSE) and hypothesis test results under both specifications. We apply our estimation methodology in [section 2.3](#) to estimate [Equation 2.7](#). For semiparametric specification, we choose the number of B-splines bases to be  $\lfloor n^{0.3} \rfloor$ .  $n$  is the number of assets in each balanced rolling window, and  $\lfloor \cdot \rfloor$  means the nearest integer. We orthogonalize these bases, and then use the Projected-PCA and restricted OLS to estimate model [Equation 2.7](#). As for the hypothesis test part, we choose threshold value to be  $\eta_n = H_n \sqrt{2 \log(PH_n)} = \lfloor n^{0.3} \rfloor \sqrt{2 \log(P \lfloor n^{0.3} \rfloor)}$ , where  $P$  is the number of characteristics, and  $n$  is the number of stock in each rolling block. For the linear specification, each characteristic only has one basis, which is itself. In terms of the hypothesis test, we use the same logic as in the semiparametric settings, and we set  $\eta_n = \sqrt{3 \log(P)}$ .

In all the estimation above, we assume that we know the real number of factors, which is three. We will discuss the situation when the number of factors is unknown in the next subsection. Mean Squared Error (MSE) is also reported to compare the fitness of models [Equation 2.7](#).

From [Table 2.1](#), under different noise levels, namely  $\sigma^2 = 1$  and  $\sigma^2 = 4$ , the semiparametric model outperforms the linear model in the following aspects:

- 1** The fitness of the semiparametric model is much better than the linear model, which can be illustrated from MSE.
- 2** The semiparametric model can enhance the power of  $S_1$  by non-zero  $S_0$ , which can not only select the correct mispricing characteristics but also avoid size distortions. As for the linear model, it is influenced by the correlated characteristics. Therefore, during certain periods we even obtain the non-invertible characteristic matrix. The linear model can also select the relevant covariates with decent probability, but it suffers from serious size distortions. In contrast, our semiparametric model with orthogonal bases can mitigate this problem to a great extent.
- 3** Because  $S_1$  can be very small and even negative, especially when the noise  $\sigma_i$  is strong, the additional component  $S_0$  is necessary to strengthen the power of  $S_1$  and to select the relevant characteristics that can explain the mispricing function.



Table 2.1 Simulation Results 1 Part I

Window	n	$\sigma^2 = 1$										$\sigma^2 = 4$														
		Linear Model					semiparametric Model					Linear Model					semiparametric Model									
		S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion %	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion %	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion %	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion %	
1	468	24.9	11.5	13.4	6.4	100%	100%	-0.5	6.2	-5.7	6	81.2%	0%	14.2	10.8	3.4	8.6	100%	100%	87.4%	-8.2	0	-8.2	8.1	0%	0%
2	894	32.8	11.6	21.2	2	100%	100%	3.4	8	-4.6	1.6	99.9%	0%	11.4	5.8	5.6	4.3	100%	100%	2.1%	-8.5	0	-8.5	3.7	0%	0%
3	1108	34.4	5.7	28.7	11.9	100%	0%	8.6	9	-0.4	11.5	100%	0%	17.1	5.7	11.4	14.1	100%	0%	0%	-7	0.7	-7.7	13.7	7.3%	0%
4	1199	-0.57	0	-0.57	10.2	0%	0%	9.2	9.1	0.1	9.5	96.8%	4.3%	-1.4	0	-1.4	12.5	0%	0%	0%	-6.1	0.06	-6.2	4.1	7%	0%
5	1333	92	19.6	72.4	2.31	100%	100%	10.6	9	1.6	2	100%	0%	28.2	6.1	22	4.5	100%	0%	0%	0.2	7.4	-7.2	4.1	88.4%	0%
6	1409	90	28.5	61.5	16	100%	100%	28.6	12.6	15.9	15.8	100%	28%	45.3	16.1	29.2	18.4	100%	73.4%	6.5%	16.3	10.9	5.4	17.5	68.4%	35.9%
7	1466	78.4	10.6	67.8	6.4	100%	100%	19.5	9	10.5	6.2	100%	0%	34.8	5.7	29.1	8.6	100%	0.02%	4.3	9	-4.7	8.4	99.9%	0%	
8	1560	133	16.8	116.2	3.3	100%	100%	20.3	10	10.3	3.2	100%	0%	45.2	6.1	39.1	5.5	100%	6.9%	4.2	10	-5.8	5.4	100%	0%	
9	1494	117.7	13.6	104.1	3.6	100%	100%	23.1	9	14.1	3.5	100%	0%	44.1	7.6	36.5	5.8	100%	32.4%	6	9	-3	5.6	100%	0.1%	
10	1292	90.7	11.5	79.2	3.7	100%	100%	16.2	9	7.2	3.6	100%	0%	39.5	9.3	30.2	5.9	100%	61.1%	6	8.9	-5.3	5.7	99.7%	0%	
11	1393	84.7	10.6	74.1	6.1	100%	85.1%	20.7	9.1	11.6	5.8	100%	1.1%	37.1	6.5	30.6	8.3	100%	12.9%	8.9	8.9	0	7.8	98.1%	1.3%	
12	1340	83.5	28	55.5	2.38	100%	100%	10.6	9	1.6	2	100%	0%	26	6.2	19.8	4.6	100%	7.1%	-1.8	5.7	-7.5	4.1	63.7	0%	
13	1285	113.8	16	97.8	1.73	100%	100%	10.6	9	1.6	1.6	100%	0%	34.5	6.6	27.9	4	100%	15.3%	-2.4	5.1	-7.5	3.7	57.1%	0%	
14	1181	88.5	12.8	75.7	4.7	100%	100%	15.8	9	6.8	4.5	100%	0%	31.2	5.9	25.3	6.9	100%	2.3%	3.7	9	-5.3	6.6	100%	0%	
15	1110	45.7	7.5	38.1	8.9	100%	30.4%	11.5	9	2.5	8.7	100%	0%	23.9	5.8	18.1	11.1	100%	0.6%	-2	4.8	-6.8	10.8	0.54%	0%	
16	1044	20.5	5.7	14.8	18.4	100%	0%	9.9	9	0.9	17.9	100%	0%	14.6	5.7	8.9	20.6	100%	0%	1.2	6.1	-4.9	20	68.1%	0.2%	
17	1125	59.4	11.5	47.9	9.2	100%	100%	13.2	9	4.2	9	100%	0%	27.2	6.2	21	11.5	100%	8.4%	2.6	8.8	-6.2	11	97.9%	0%	
18	2192	NA	NA	NA	NA	NA	NA	23.2	11	12.2	4.3	100%	0%	NA	NA	NA	NA	NA	NA	NA	6.7	11	-4.3	6.4	100%	0%
19	2236	56.1	11.5	44.6	5.8	100%	100%	17.8	11	6.8	5.2	100%	0%	28.3	6.3	22	8	100%	20.3%	4.3	11	-6.7	7.4	100%	0%	
20	2273	43.3	5.7	37.6	3.8	100%	0%	22.4	11	11.4	3.2	100%	0%	22.4	5.7	16.7	6.1	100%	0%	5	10.2	-5.2	5.4	92.6%	0%	
21	2235	59.8	11.8	48	2.7	100%	100%	20.2	11	9.2	2	100%	0%	25	7.3	17.7	4.9	100%	28.2%	4.6	11	-6.4	4.2	100%	0%	
22	2270	40.2	11.5	28.7	2.78	100%	99.5%	17.2	11.6	5.6	2.1	100%	0%	17.1	5.9	11.2	5	100%	3.5%	-6	0.1	-6.1	4.2	1.1%	0%	
23	2405	41.4	8.9	32.5	4.1	100%	54.2%	16.3	11	5.3	3.3	100%	0%	18.7	5.8	12.9	6.3	100%	7.1%	-3.3	3	-6.3	5.5	27.3%	0%	
24	2376	19	9.7	9.3	1.8	100%	69.9%	23.1	11	12.1	1	100%	0%	7.5	5.7	1.8	4	100%	0%	5.6	11	-5.4	3.2	100%	0%	
25	2323	15.9	9.5	6.4	3.5	66.7%	98.6%	20.6	11	9.6	2.7	100%	0%	1.1	0	1.1	5.8	0%	0%	5.3	11	-5.7	4.9	100%	0%	
26	2344	NA	NA	NA	NA	NA	NA	24.9	12.9	12	3.3	100%	17.1%	NA	NA	NA	NA	NA	NA	NA	6.5	11	-4.5	5.4	100%	0%
27	2434	NA	NA	NA	NA	NA	NA	27.3	11	16.3	1.2	100%	0%	NA	NA	NA	NA	NA	NA	NA	6.9	11	-4.1	3.4	100%	0%
28	2548	0.9	0	0.9	4.2	0%	0%	26.2	11	15.2	3.3	100%	0%	-1.3	0	-1.3	6.5	0%	0%	6.9	11	-4.1	5.5	100%	0%	
29	2741	10.3	5.7	4.5	4.2	100%	0%	58.2	11.1	47.1	3.4	100%	1.3%	6.6	5.7	0.9	6.4	100%	1.3%	17.6	11	6.6	5.5	100%	0%	
30	2928	5.6	4.6	1	7.1	80.4%	0%	59.2	11.8	47.4	6.3	100%	7.8%	-0.4	0.1	-0.5	9.3	2.5%	21.6%	18.8	11	7.8	8.5	100%	0.3%	
31	2894	13.4	5.7	7.7	6.4	100%	0%	61	13.4	47.6	5.7	100%	0%	8.1	5.7	2.3	8.6	100%	0%	17.7	11	6.7	7.8	100%	0.2%	
32	2905	23.1	11.5	11.6	5.9	100%	100%	33.2	11.3	21.9	5.2	100%	3.6%	12.9	8.5	4.4	8.1	100%	48.2%	9.8	11	-1.2	7.4	100%	0%	
33	2804	9.8	5.7	4.1	9.6	100%	0%	42.7	18.5	24.2	8.9	100%	68.5%	7.3	5.7	1.6	11.9	100%	0%	9.7	11	-1.3	11.2	100%	0%	
34	2570	6.9	5.7	1.2	2.2	99.7%	0%	37.3	12.2	25.1	21.2	100%	10.4%	2	1.9	0.1	2.4	34.4%	0%	12.7	11	1.7	23.3	100%	0.2%	
35	2516	8.3	5.7	2.6	7.9	100%	0%	41.3	11	30.3	7.2	100%	0.4%	5.1	5.02	0.28	10.1	87.3%	0%	12.9	11	1.9	9.4	100%	0%	
36	2491	10.7	5.7	4.9	2.1	100%	0%	41.3	11	30.3	1.4	100%	0.4%	0.5	0.25	0.25	4.4	4.5%	0%	12.4	11	1.4	3.6	100%	0%	
37	2402	14.1	5.7	8.4	5.6	100%	0%	26.5	11.2	15.3	4.9	100%	2.2%	8.8	5.7	3.1	7.9	100%	2.2%	8.8	7.9	11	-3.1	7.1	100%	0%
38	2326	19.7	9.6	10.1	3	100%	66.8%	28.9	11.3	17.6	2.3	100%	2.1%	8.1	5.8	2.3	5.3	100%	0.3%	8.7	11	-2.3	4.4	99.9%	0.1%	
39	2241	17	5.7	16.1	2.9	100%	0.2%	11	11	0	1.7	100%	0%	9.1	5.8	2.3	5.3	100%	0.3%	-7.5	0.1	-7.6	4	1.1%	0%	
40	2178	21.8	5.7	16.1	2.9	100%	0%	9.5	11	-1.5	2.2	100%	0.3%	12.2	5.7	6.5	5.2	100%	0%	-8.1	0	-8.1	4.4	0%	0%	

Table 2.2 Simulation Results 1 Part2

Window	n	$\sigma^2 = 1$										$\sigma^2 = 4$													
		Linear Model					semparametric Model					Linear Model					semparametric Model								
		S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion%	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion%	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion%	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Selected %	Distortion%
41	2113	24.1	6.1	18	4.7	100%	7.5%	7.8	10	-2.2	3.9	100%	0%	13.9	5.7	8.2	6.9	100%	0%	-8.1	0	-8.1	6.1	0%	0%
42	2023	18.4	5.7	12.7	6.8	100%	0%	11.3	10	1.3	6	100%	0%	10.8	5.8	5.1	9	100%	0%	-7.1	0.3	7.4	8.2	2.7%	0%
43	2007	18.8	5.7	13.1	4.9	100%	0%	9.1	10	-0.9	4.1	100%	0%	10.5	5.7	4.8	7.1	100%	0%	-8.3	0	-8.3	6.3	0%	0%
44	1924	16.6	5.8	10.8	8.18	100%	0.2%	13.6	10.8	2.8	7.5	100%	8%	11.2	5.8	5.4	10.4	100%	0.3%	-3.5	2.7	-6.2	9.7	26.3%	0.2%
45	1990	27.5	5.7	21.8	2.1	100%	0%	8.1	10	-1.9	1.4	100%	0%	13.3	5.7	7.5	4.4	100%	0%	-8	0	-8	3.6	0%	0%
46	1937	20.3	5.8	14.5	5.4	100%	0.9%	19.7	11.8	7.9	4.7	100%	18%	12.6	5.9	6.7	7.6	100%	3%	8	11.2	-3.2	6.8	100%	12.3%
47	1909	13.2	5.7	7.5	5.2	100%	0%	14.2	10.4	3.8	4.5	100%	3.5%	8.8	5.7	3.1	7.4	100%	0%	2.7	8.4	-5.7	6.7	84.9%	0%
48	1872	21.8	5.7	16.1	2.7	100%	0%	11.4	10	1.4	2	100%	0%	11.1	5.8	5.3	4.9	100%	0%	-6.8	0.6	-7.4	4.2	5.7%	0%
49	1841	16.3	5.7	10.5	2.1	100%	0%	8.7	10	-1.3	1.4	100%	0.1%	8.1	5.7	2.4	4.4	100%	0%	-8.4	0	-8.4	3.6	0%	0%
50	1826	11	5.7	5.3	4.3	100%	0%	12.6	10.6	2	3.5	100%	3.5%	6.5	5.7	0.8	6.6	99.7%	0.3%	-6.9	0	-6.9	5.7	0%	0%

This table documents results under the characteristics-based beta and alpha of Fama-French 3 factors model. To mimic the empirical study, we simulated 30 12-month rolling windows, and each window is repeated for 1000 times. Each column summarises the mean value of 1000 estimations and test results. S<sub>1</sub> is the conventional Wald test while S<sub>0</sub> is the power-strengthened component. This table also compares the performance of both semiparametric and linear models under different noise levels,  $\sigma^2 = 1$  and  $\sigma^2 = 4$ . NA results are caused by non-invertible characteristic matrices. "Selected" means the percentage of selecting the relevant characteristic in the mispricing function in 1000 experiments. Similarly, "distortion" represents the percentage of wrongly selecting irrelevant characteristics in 1000 repetitions.

### 2.7.3 Robustness Under Stronger Noise

In [Table 2.1](#), we set two different noise levels of random shocks, namely  $\sigma^2 = 1$  and  $\sigma^2 = 4$ . Although  $\sigma^2 = 1$  is closer to the empirical data, we conduct this comparison to show the robustness of our methods. When the noise level becomes three times bigger, the accuracy of power-enhanced tests gets much lower for certain windows. However, there are no size distortions under comparatively high noise level recalling that all the components of our simulation model are rescaled to have unit variance. Another fact is that the stronger noise does deteriorate the power of conventional Wald tests, leading to an even smaller value of  $S_1$ , which can be mitigated through adding  $S_0$ .

Therefore, we conclude that our methods are robust to a higher noise level regarding no size distortions. However, the accuracy of selecting relevant components and the role of enhancing the power of hypothesis tests will be influenced negatively.

### 2.7.4 Number of Factors

In the empirical study, the number of factors is unknown. Therefore, in this subsection we will study whether our methodology is robust to a various numbers of factors considered.

We simulate according to another data generation process:

$$y_{it} = h(X_i) + \sum_{j=1}^5 g_j(\mathbf{X}_{ij})f_{jt} + \varepsilon_{it}, \quad (2.8)$$

similarly,  $y_{it}$  is the generated stock return;  $h(X_i)$  is the mispricing function consisting of a nonlinear characteristic function of  $X_i$ , to mimic the sparse structure of the mispricing function;  $g_j(\mathbf{X}_{ij})$  is the  $j^{\text{th}}$  characteristics-based factor loading, which has an additive semiparametric structure;  $X_{ij}$  is a subset consisting of four characteristics;  $f_{jt}$  is the  $j$  Fama-French 5-factor returns at time  $t$ ;  $\varepsilon_{it}$  is the idiosyncratic shock, generated from  $N(0, \sigma^2)$ . Moreover, we generate characteristic functions as:

$$h(X_i) = \sin X_i,$$

$$g_j(\mathbf{X}_j) = X_{ij1}^2 + (3X_{ij2}^3 - 2X_{ij2}^2) + (3X_{ij3}^3 - 2X_{ij3}) + X_{ij4}^2,$$

where  $X_{ijl}$  is a randomly picked characteristic without replacement from the data in empirical study with  $j = 1, \dots, 5$ ,  $l = 1, \dots, 4$ . Furthermore, all  $h(X_i)$  and  $g_j(\mathbf{X}_{ij})$  are rescaled to have zero mean and unit variance.

Given the above data generation process, together with the data generation process, we test the influence of over and under-estimated number of factors. We choose the number of factors to be either three or five, and compare the results in [Table 2.3](#). The first category column is the scenario of over estimating the number of factors. We simulate the data generation process using the Fama-French three factors model, but estimate the number of factors to be five. However, this does not cause any serious problems. For some rolling blocks, the probability of mistakenly selected irrelevant characteristics is slightly higher under over estimating the number of factors. Moreover, over estimating the number of factors can increase the model fitting marginally. Therefore, we conclude that over estimating the number of factors does not cause severe size distortion using our methods.

On the other hand, under estimating the number of factors can lead to misleading test results. We can conclude this from the last column where we estimate the number of factors to be three in a five-factor model. Compared with the correct specified model, under estimating causes not only higher MSE, but also higher distortions, which means it is more likely to select irrelevant characteristics. Therefore, in the empirical study we prefer the five-factor model rather than the three-factor model.





## 2.8 Empirical Study

### 2.8.1 Data

We use monthly stock returns from CRSP and firms' characteristics from Compustat, ranged from 1965 to 2017. We construct 33 characteristics following the methods of [Freyberger et al. \(2020a\)](#). Details of these characteristics can be found in the Appendix. We use characteristics from fiscal year  $t - 1$  to explain stock returns between July of year  $t$  to June of year  $t + 1$ . After adjusting the dates from the balance sheet data, we merge two data sets from CRSP and Compustat. We require all firms included in our analysis to have at least three years of characteristics data in Compustat.

Data is modified with regards to the following aspects:

- 1 Delisting is quite common for CRSP data. We use the way of [Hou et al. \(2015\)](#) to correct the returns of all delisted assets before 2018. Detailed methods can be found in their Appendix.
- 2 Projected-PCA works well, even under small  $T$  circumstances. Thus, we choose the width of our window to be 12 months. Another reason for the short window width is that we assume that mispricing functions are time-invariant in each window. One of the limitations of Projected-PCA is that it can only be used for a balanced panel, which means the number of stock will vary when we applied one-year rolling windows to obtain a short-time balanced panel. Meanwhile, we take monthly updated characteristics' mean values of 12 months as fixed characteristic values in each window. We also use the rolling window method to detect peer groups of arbitrage characteristics.
- 3 B-splines are based on each time-invariant characteristic of  $n$  firms, which are not delisted in each window.
- 4 Rolling windows are moving at a 12-month step from Jul. 1967 to Jun. 2017 without overlapping. The first 24 months returns are not included as they do not have corresponding characteristics.
- 5 Excess returns are obtained by the difference between monthly stock returns and monthly Fama-French risk-free returns.

### 2.8.2 Estimation

We construct B-splines bases based on evenly distributed knots, and the degree of each basis is three. We choose  $H_n = \lfloor n^{0.3} \rfloor$ , and  $n$  is the number of stocks. To get a relatively large balanced panel in each window, some characteristics with too many missing values are eliminated. Therefore, only 33 characteristics are left. Firms kept in balanced panels in our dataset range from 468 to 2928, which means that both  $n$  and  $\hat{\mathbf{A}} \in \mathbb{R}^{PH_n}$  are diverging. Large  $n$  can satisfy asymptotic requirements. These facts emphasize the necessity of introducing a power-enhanced component into the hypothesis test. Before the next step, we use time-demeaning matrix  $\mathbf{D}_T$  to demean excess return matrix in each window.

Next, we project the time-demeaned monthly excess return matrix  $\tilde{\mathbf{Y}}$  onto the B-splines space spanned by characteristics bases  $\Phi(\mathbf{X})$ , and then we collect the fitted values  $\hat{\mathbf{Y}}$ . We apply Principal Component Analysis on  $\frac{1}{n}\hat{\mathbf{Y}}^\top\hat{\mathbf{Y}}$ , and attain the first five eigenvectors corresponding to the first five biggest eigenvalues as the estimates of unobservable factors  $\mathbf{F}$ . We choose the number of factors to be five according to simulation results.

Then, we estimate the factor loading matrix by:

$$\hat{\mathbf{G}}(\mathbf{X}) = \hat{\mathbf{Y}}\hat{\mathbf{F}}(\hat{\mathbf{F}}^\top\hat{\mathbf{F}})^{-1}.$$

Moreover, we use equality-constrained OLS to estimate the mispricing function. We project excess monthly return matrix on the characteristic space  $\Phi(\mathbf{X})$  that is orthogonal to factor loading matrix  $\hat{\mathbf{G}}(\mathbf{X})$ .

Another goal of this paper is to conduct a power-enhanced test on the mispricing function. Therefore, our final step is to estimate the covariance matrix  $\Sigma$  of  $\hat{\mathbf{A}}$ .

### 2.8.3 Power-enhanced Hypothesis Tests

In this section, we conduct a power-enhanced test in each rolling block. Firstly, we set threshold value for each window,  $\eta_n = H_n\sqrt{2\log(PH_n)}$ , where  $H_n$  is the number of bases for each characteristic whereas  $P$  is the number of total characteristics in each window, with  $P = 33$ .  $\eta_n$  is data-driven critical value and it diverges as the number of firms increases. We use indicator function  $\mathbf{I}(\sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n)$  with critical value  $\eta_n = H_n\sqrt{2\log(PH_n)}$  to achieve three goals.



- 1 This indicator function select the most relevant characteristics that can explain the variation of the mispricing function. Results of last column in [Table 2.5](#) are characteristics selected in  $\hat{\mathcal{M}} = \{\mathbf{X}_p \in \mathbf{X} : \sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n, \quad p = 1, 2, \dots, P\}$ .
- 2 It contributes to the test statistics  $S$  by adding a diverging power-enhanced component  $S_0$ . As  $T = 12$  is small in this empirical study, we assume the homoskedasticity of  $\varepsilon_{it} + \gamma_i$ . We also specify an over-shrunk covariance matrix by setting off-diagonal elements to be zeros.
- 3 It avoids size-distortion by the conservative critical value  $\eta_n$ .

The diagonal elements of  $\hat{\Sigma}$  are estimated variances of mispricing coefficients. These elements can be substituted into the indicator function  $\mathbf{I}(\sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n)$ , where  $\hat{\sigma}_{ph}$  is the  $ph^{th}$  diagonal element of  $\hat{\Sigma}$ .

Finally, the new statistics  $S$  can be calculated as:

$$S = S_0 + S_1,$$

$$S_0 = H_n \sum_{p=1}^P \mathbf{I}\left(\sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n\right), \quad S_1 = \frac{\hat{\mathbf{A}}^\top \hat{\Sigma}^{-1} \hat{\mathbf{A}} - PH_n}{\sqrt{2PH_n}}.$$

#### 2.8.4 Test Results

This section presents the empirical results. Details can be found in [Table 2.5](#), which lists the results of 50 rolling windows from Jul.1967 to Jun.2017. Generally, the number of firms included in the 12-month rolling block is increasing period by period. The number of our characteristic B-splines bases is a function of the number of firms  $n$  in each block. Therefore, the dimension of the mispricing coefficient vector  $\mathbf{A} \in \mathcal{R}^{PH_n}$  is also diverging. This verifies the necessity of using a power-enhanced component  $S_0$ .

Recalling that  $S|H_0 \rightarrow^d N(0, 1)$ , some of the test statistics  $S$  is big enough to reject the null hypothesis. However, for some testing windows, there are no strong signals showing the existence of characteristics-based mispricing functions after subtracting systematic effects. Most  $S_1$  values are small and even negative, which may be caused by the sparsity structure of the mispricing function or/and the low power problems due to diverging dimension of mispricing coefficients.

The power-enhanced component  $S_0$  works well in the empirical study. It selects the most important explanatory characteristics and strengthens the power of  $S_1$ , mitigating the low power problem.

Apart from contributing to the power of tests, the indicator function in the power-enhanced component can also screen out the most relevant characteristics, which are concluded as "Characteristics Selected" in [Table 2.5](#).

Some empirical findings are worth discussing. Although short-term cumulative returns like  $r_{2\_1}$  are always selected, we cannot take this as the evidence of arbitrage opportunities since we construct  $r_{2\_1}$  as the time-invariant average of all  $r_{2\_1}$  in the same rolling window, which contains much overlapping information of monthly excess returns. However, this is not the case for long-term and mid-term cumulative returns like  $r_{12\_2}$ ,  $r_{12\_7}$  and  $r_{6\_2}$ , because these average returns include a lot of information from another rolling window.

Apart from the cumulative returns, some other characteristics contribute to the arbitrage opportunities as well. PCM (Price to Cost Margin) appears twice. From [Figure B.2](#), we find that the PCM mispricing curve is nonlinear and generally decreasing as the value of PCM increases. ROA (Return-On-Asset) also plays a role during 1988-1989. It behaves like a parabola with fluctuations near zero in [Figure B.3](#). As for Lev (ratio of long-term debt and debt in the current liabilities), it is decreasing for Lev 0 and increasing afterward as in [Figure B.7](#). In [Figure B.8](#), IPM (pre-tax Profit Margin) function behaves like a "V" shape with the turning point zero during 2004-2005. DelGmSale (difference in the percentage in Gross margin and the percentage change in Sales) experiences a bump at zero during 2015-2016 in [Figure B.9](#). C2D (cash flow to price) curve behaves like "V" around the zero in 2016-2017, (see [Figure B.10](#)). All characteristics curves in the above figures are standardized as uniformly distributed characteristics in the interval  $[-100, 100]$ . This is for presentation purposes only since most characteristics are unevenly distributed.

Another finding is the persistence of some arbitrage characteristics. Arbitrage characteristics can be persistent for two years once appeared, such as BEME (ratio of the book value of equity and market value of equity) in [Figure B.4](#). Some persistent arbitrage characteristics even have similar shapes of mispricing functions in different rolling windows, such as AT (Total asset) in [Figure B.6](#) and LME (total market capitalization of the previous month) in [Figure B.5](#).

Table 2.5 Empirical Study Results

Time period	n	S	S <sub>0</sub>	S <sub>1</sub>	MSE	Characteristics Selected
Jul.1967-Jun.1968	468	-9.6	0	-9.6	0.005	NONE
Jul.1968-Jun.1969	951	-0.45	8	-8.45	0.004	$r_{2\_1}$
Jul.1969-Jun.1970	1108	1.7	9	-7.3	0.005	$r_{2\_1}$
Jul.1970-Jun.1971	1199	-8.7	0	-8.7	0.006	NONE
Jul.1971-Jun.1972	1333	-10	0	-10	0.004	NONE
Jul.1972-Jun.1973	1409	12.7	18	-5.3	0.005	$r_{12\_2}, r_{6\_2}$
Jul.1973-Jun.1974	1466	2.1	9	-6.9	0.005	$r_{2\_1}$
Jul.1974-Jun.1975	1560	-10.7	0	-10.7	0.01	NONE
Jul.1975-Jun.1976	1494	0.1	9	8.9	0.05	$r_{2\_1}$
Jul.1976-Jun.1977	1292	0.1	9	-9	0.004	$r_{2\_1}$
Jul.1977-Jun.1978	1393	-9.4	0	-9.4	0.005	NONE
Jul.1978-Jun.1979	1340	8.6	18	-9.4	0.005	$r_{2\_1}, r_{12\_7}$
Jul.1979-Jun.1980	1285	1	9	-8	0.005	$r_{2\_1}$
Jul.1980-Jun.1981	1181	9.7	18	-8.2	0.006	$r_{12\_7}, r_{12\_2}$
Jul.1981-Jun.1982	1110	1.2	9	-7.8	0.01	$r_{2\_1}$
Jul.1982-Jun.1983	1044	33.1	36	-3	0.01	$r_{12\_2}, r_{12\_7}, r_{6\_2}, r_{2\_1}$
Jul.1983-Jun.1984	1125	-0.9	9	-9.9	0.006	$r_{2\_1}$
Jul.1984-Jun.1985	2192	-0.2	11	-11.2	0.01	$r_{2\_1}$
Jul.1985-Jun.1986	2236	13.1	22	-8.94	0.01	$r_{12\_7}, r_{12\_2}$
Jul.1986-Jun.1987	2273	1.7	11	-9.3	0.01	PCM
Jul.1987-Jun.1988	2235	0.9	11	-10.1	0.01	$r_{2\_1}$
Jul.1988-Jun.1989	2270	1.2	11	-9.8	0.01	ROA
Jul.1989-Jun.1990	2405	-0.1	11	-11.1	0.01	$r_{2\_1}$
Jul.1990-Jun.1991	2376	1.1	11	-9.9	0.02	$r_{2\_1}$
Jul.1991-Jun.1992	2323	2.1	11	-8.9	0.02	$r_{2\_1}$
Jul.1992-Jun.1993	2344	12.2	22	-9.8	0.02	$r_{12\_7}, r_{12\_2}$
Jul.1993-Jun.1994	2434	0.4	11	-10.6	0.01	$r_{2\_1}$
Jul.1994-Jun.1995	2548	2.4	11	-8.6	0.01	$r_{2\_1}$
Jul.1995-Jun.1996	2741	14.1	22	-7.9	0.02	BEME, $r_{2\_1}$
Jul.1996-Jun.1997	2928	18.1	22	-3.9	0.01	BEME, $r_{2\_1}$
Jul.1997-Jun.1998	2894	26.5	33	-6.5	0.02	$r_{2\_1}, r_{12\_7}, r_{12\_2}$
Jul.1998-Jun.1999	2905	24.6	33	-8.4	0.02	AT, LME, $r_{2\_1}$
Jul.1999-Jun.2000	2804	13.8	22	-8.2	0.03	$r_{2\_1}, r_{12\_7}$
Jul.2000-Jun.2001	2570	37.7	44	-6.3	0.02	AT, LME, $r_{2\_1}, r_{6\_2}$
Jul.2001-Jun.2002	2516	1.3	11	-9.7	0.02	$r_{2\_1}$
Jul.2002-Jun.2003	2491	15	22	-7	0.02	Lev, $r_{2\_1}$
Jul.2003-Jun.2004	2402	3.9	11	-7.1	0.01	$r_{2\_1}$
Jul.2004-Jun.2005	2326	1.8	11	-9.2	0.01	IPM
Jul.2005-Jun.2006	2241	2.5	11	-8.5	0.01	$r_{2\_1}$
Jul.2006-Jun.2007	2178	1.5	11	-9.5	0.01	$r_{2\_1}$
Jul.2007-Jun.2008	2113	12.6	20	-7.4	0.01	$r_{12\_2}, r_{2\_1}$
Jul.2008-Jun.2009	2023	1.7	10	-8.3	0.02	$r_{2\_1}$
Jul.2009-Jun.2010	2007	1	10	-9	0.01	$r_{2\_1}$
Jul.2010-Jun.2011	1924	13.6	20	-6.4	0.01	$r_{2\_1}$
Jul.2011-Jun.2012	1990	2.5	10	-7.5	0.01	$r_{2\_1}$
Jul.2012-Jun.2013	1937	23.7	30	-6.3	0.01	$r_{2\_1}, r_{12\_7}, r_{12\_2}$
Jul.2013-Jun.2014	1909	2.3	10	-7.7	0.01	$r_{2\_1}$
Jul.2014-Jun.2015	1872	5.5	10	-4.5	0.01	$r_{2\_1}$
Jul.2015-Jun.2016	1841	12.4	20	-7.6	0.01	DelGmSale, $r_{2\_1}$
Jul.2016-Jun.2017	1826	26.1	30	-3.9	0.01	C2D, PCM, $r_{12\_7}$

This table summarizes the empirical results, where n represents the number of stock in this rolling window.

### 2.8.5 Dynamic Peer Groups of Arbitrage Characteristics

In this section, we illustrate that there are distinguishable peer groups of the same arbitrage characteristic resulting in similar mispricing returns. We apply the methods in [section 2.5](#) and take two rolling windows, namely, Jul.1986- Jun.1987 and Jul.2004-Jun.2005 as demonstrative examples.

In the rolling window Jul.1986-Jun.1987, PCM is selected as the only arbitrage characteristic that can explain arbitrage returns. We reveal that similar characteristic-based arbitrage returns are determined by distinguishable groups of the characteristic PCM. We first divide arbitrage returns  $\ddot{y}_{it}$  into different return groups. And then, we detect whether there are some clustering structures within groups of the highest and the lowest characteristic-based arbitrage returns, respectively. As we have 2326 assets, for the visualization purpose, we set the threshold value of the K-means method to be relatively small to have as many as ten groups.

Table 2.6 First layer 1986-1987 (clusterings of  $\ddot{y}_{it}$  )

Group number	Group centroid	Group size
1	0.0059	435
2	0.1205	26
3	-0.0082	428
4	0.0399	189
5	0.0697	71
6	-0.1018	29
7	-0.0617	110
8	-0.0390	250
9	-0.0225	349
10	0.0208	386

In [Table 2.6](#), group 2 has the largest positive average return while group 6 has the worst. Next, we detect the clusterings of characteristic "PCM" within each group individually, which is the second layer in [section 2.5](#).

Table 2.7 Second layer 1986-1987 (clusterings of characteristic PCM )

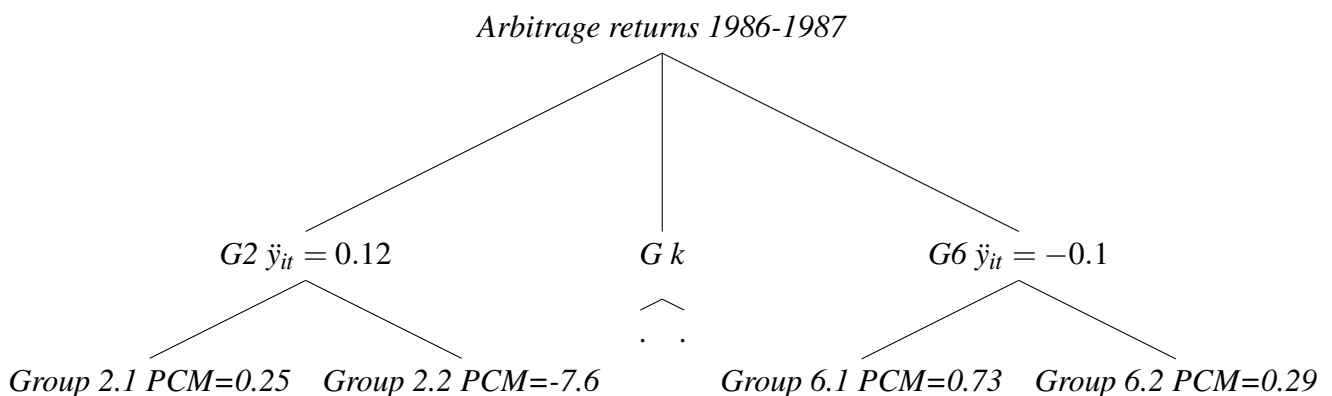
Group number	Centroids of Arbitrage returns	Centroids of PCM	Group size
2.1	0.1211	0.2452	25
2.2	0.1039	-7.630	1

In [Table 2.7](#), there are two clusterings of PCM that provide the highest positive characteristic-based arbitrage returns. Group 2.2, which has an extreme negative PCM value but a high characteristic-based arbitrage return, is an outlier. Members in group 2.1 with excellent arbitrage performance have positive and small PCM values.

Table 2.8 Second layer 1986-1987 (clusterings of characteristic PCM )

Group number	Centroids of Arbitrage returns	Centroids of PCM	Group size
6.1	-0.1085	0.728	9
6.2	-0.0989	0.288	20

[Table 2.8](#) gives groups of PCM in group 6. Members of this group are divided into two clusterings. Group 6.1 has a relatively large PCM value, while group 6.2 has a smaller PCM, which is close to that in group 2.1 with the highest arbitrage return. This is an evident illustration of the nonlinear structure of  $h(\mathbf{X})$  in this window. The structure of characteristic-based arbitrage returns during Jul.1986- Jun.1987 is:



The classification can be found at [Figure B.11](#), where assets are labeled by their "PERMNO," and both axes are rescaled.

Another example is the characteristic-based arbitrage return  $\ddot{y}_{it}$  during the year 2004-2005. The power-enhanced test selects characteristic "IPM" as the only explanatory variable.

We apply the Hierarchical K-means method. The results of the first layer classification can be found in [Table 2.9](#). There are ten groups in total according to the similarity of characteristic-based arbitrage returns. Next, we pick two groups with the highest and the lowest returns, respectively, to give clusterings of "IPM" in these two groups.

Table 2.9 First layer 2004-2005 (clusterings of  $\ddot{y}_{it}$  )

Group number	Group centroid	Group size
1	0.0421	276
2	0.0059	459
3	0.1537	26
4	-0.024	367
5	0.0659	166
6	0.023	387
7	0.0999	120
8	-0.0758	67
9	-0.0437	244
10	-0.0082	436

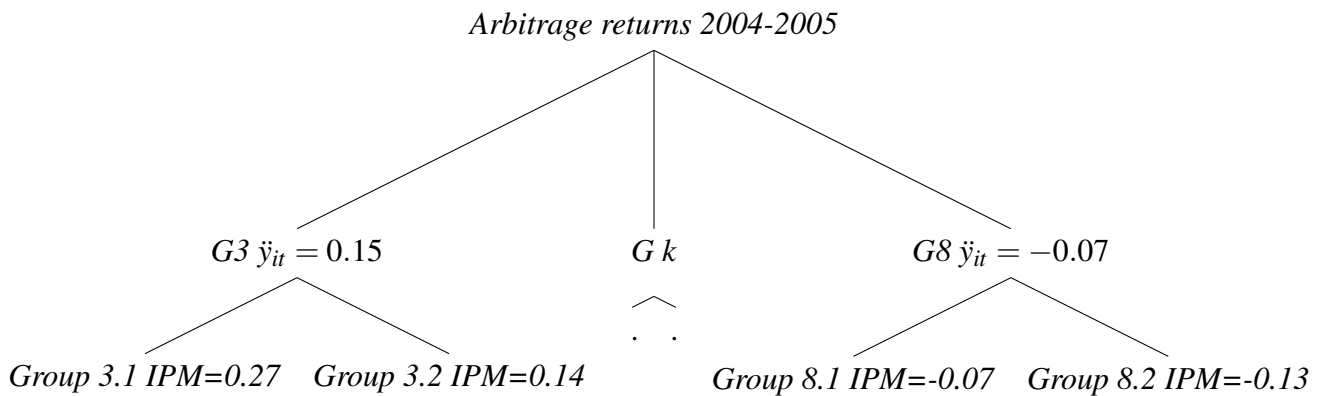
Similarly, we show classification results in [Table 2.10](#) and [Table 2.11](#). Positive IPM values give higher characteristic-based arbitrage returns. On the contrary, when IPM is close to zero or negative, the characteristic-based arbitrage returns fall into the lowest group (group 8).

Table 2.10 Second layer 2004-2005 (clusterings of characteristic IPM )

Group number	Centroids of Arbitrage returns	Centroids of PCM	Group size
3.1	0.1681	0.266	5
3.2	0.1502	0.143	21

Table 2.11 Second layer 2004-2005 (clusterings of characteristic IPM )

Group number	Centroids of Arbitrage returns	Centroids of PCM	Group size
8.1	-0.0713	-0.07	10
8.2	-0.1016	-0.134	57



The plots of the IPM can be found at [Figure B.12](#), where the axes are rescaled, and assets are labeled by their "PERMNO" code with five digits.

Finally, it is obvious that peer groups of arbitrage characteristics are dynamic in two aspects. Firstly, the selected arbitrage characteristics are time-varying. Although some of the arbitrage characteristics can show up for more than one block once they appeared, no arbitrage characteristic can be substantially persistent. Secondly, as in [Figure B.4](#), the same arbitrage characteristic can have different functional forms in various rolling windows. However, the patterns of some characteristics show persistence in different time periods, such as AT in [Figure B.6](#) and LME in [Figure B.5](#). In a word, under the flexible semiparametric setting, methods for constructing arbitrage portfolio in [Kim et al. \(2019\)](#) are inapplicable, although the characteristics-based mispricing function is significant for certain time periods.

## 2.9 Conclusion

We proposed a semiparametric characteristics-based factor asset pricing model, with a focus on the existence and the structure of the mispricing function. Both unknown characteristics-based factor loadings and the mispricing component are approximated by B-splines sieve. The model is estimated by both Project-PCA and equality-constrained OLS. We also develop a power-enhanced test to investigate whether there are mispricing characteristics, orthogonal to the main systematic factors. This is necessary because when the B-splines coefficients of the mispricing function are diverging, the conventional Wald test has very low power. The traditional Wald test is strengthened by a screening process, which adds significant components to the original statistics. Our proposed methods work well for both simulations and the US stock market. Empirically, we find distinct clusterings of the same characteristics resulting in similar arbitrage returns, forming "peer groups." The mispricing functions are time-varying and mostly insignificant under our setting.





## **Chapter 3**

# **A Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection**



## Abstract

This paper develops a two-step semiparametric methodology for portfolio weight selection for characteristics-based factor-tilt and factor-timing investment strategies. We build upon the expected utility maximization framework of [Brandt \(1999\)](#) and [Aït-sahalia and Brandt \(2001\)](#). We assume that asset returns obey a characteristics-based factor model with time-varying factor risk premia as in [Ge et al. \(2020\)](#). We prove under our return-generating assumptions that an approximately optimal portfolio can be established using a two-step procedure in a market with a large number of assets. The first step finds optimal factor-mimicking sub-portfolios using a quadratic objective function over linear combinations of characteristics-based factor loadings. The second step dynamically combines these factor-mimicking sub-portfolios based on a time-varying signal, using the investor's expected utility as the objective function. We develop and implement a two-stage semiparametric estimator. We apply it to CRSP (Center for Research in Security Prices) and FRED (Federal Reserve Economic Data) data and find excellent in-sample and out-sample performance consistent with investors' risk aversion levels.

KEYWORDS: Portfolio management; Prediction; Single index; GMM;

JEL CLASSIFICATION: C14; G11.

### 3.1 Introduction

The traditional portfolio choice model proceeds by estimating the parameters of an asset return distribution and then finding the portfolio that maximizes expected payoffs for a given risk level, such as the optimal mean-variance portfolio choice model proposed by [Markowitz et al. \(1952\)](#). This approach can produce biases in portfolio weights since the portfolio selection process ignores the estimation error in the empirically-derived return distribution parameters. Furthermore, as the number of assets increases, the estimation of the high-dimensional covariance matrix becomes intractable. Some notable methods have been proposed to solve this issue, such as linear ([Ledoit and Wolf \(2004\)](#)) and nonlinear shrinkage ([Ledoit and Wolf \(2017\)](#)) of the target covariance matrix or selecting main elements by threshold ([Fan et al. \(2013\)](#)). However, these approaches may cause information loss and lead to unsatisfactory results, as illustrate by [Ao et al. \(2019\)](#). At the same time, [Ao et al. \(2019\)](#) studied a method called MAXSER, which is a sparse regression that sets the optimal Sharpe-ratio as the regressand. Their method also requires a sparsity assumption and can be problematic when the number of assets  $n$  is large. Meanwhile, all aforementioned papers ignore the importance of predictive variables, which have been documented by many researchers, such as [Fama and French \(1989\)](#), who analyzed the forecasting ability of dividend yield, default spread, and term spread on asset returns, as well as [Keim and Stambaugh \(1986\)](#), [Campbell and Shiller \(1988\)](#), [Hodrick \(1992\)](#), [Chen et al. \(2016\)](#), [Gu et al. \(2020\)](#) and [Chen et al. \(2020\)](#) among others. The goal of this paper is to construct an optimal two-step portfolio that takes advantage of both a large number of assets and dynamic predictors.

[Brandt \(1999\)](#) used nonparametric tools to directly estimate the portfolio weights that maximize the expected utility of the observed data, without first estimating the return distribution. He estimated the dynamic portfolio weights of the assets in a two-asset model as a nonparametric function of the uni-variate time-series predictor of the future excess returns of the risky assets. [Aït-sahalia and Brandt \(2001\)](#) replaced the uni-variate time series predictor with an index-based set of predictors: the time-varying portfolio weights in a three-asset model were assumed to be a nonlinear function of a linear fixed combination of a vector of predictive variables. However, the number of assets included in their portfolio was quite limited.

[Brandt et al. \(2009\)](#) developed a characteristic-based model for portfolio selection with a large cross-section of assets. They assumed that optimal portfolio weights were linearly related to a small set of observable characteristics, such as book-to-market ratio, momentum, and market capitalization. They found the linear coefficients that maximized expected utility under this assumption.

In this paper, we develop a new semiparametric model of portfolio selection, which combines the advantages of a large cross-section of assets and dynamic predictive variables. This is achieved by a characteristics-based asset pricing factor model. We generalize the methodologies in the papers mentioned above since we do not impose the assumption that optimal weights are linear in the characteristics. Furthermore, the firm-specific characteristics included in our model can be significantly broadened. There are 33 characteristics in our empirical study, which provide more potential abnormal return opportunities. Also, as in [Aït-Sahalia and Brandt \(2001\)](#), we also allow information-based dynamically-varying portfolio allocation based on a single-index function of predictors. We replace weighting across asset classes in [Aït-Sahalia and Brandt \(2001\)](#) with weighting across our optimally-constructed characteristics-based sub-portfolios.

We estimate the model using a new, two-stage semiparametric procedure. The first step involves the estimation of the factor-mimicking sub-portfolios, which is a high-dimensional estimation problem since the number of assets is diverging. Still, the objective function is quadratic, allowing us to solve it using semiparametric techniques. That step compacts those assets into several sub-portfolios rather than discarding some of them and reduces dimensionality, simplifying the next step. The second step maximizes the dynamic expected utility of a risk-free asset and those sub-portfolios conditional on a set of predictors, similar to [Aït-Sahalia and Brandt \(2001\)](#). Our two-step statistical methodology accounts fully for the estimation error in both semiparametric steps, and we show that it approximates the intractable single-stage, asset-by-asset portfolio weight estimation problem in a well-defined sense.

Our model is not entirely general: we do not allow individual asset selection in response to asset-specific valuation information. We essentially allow for factor-tilt strategies, which means weighting securities according to their factor exposure in response to the associated factor risk premia, and factor timing, which means dynamically varying factor-tilt strategies, accounting for predictability in factor risk premia, but not individual asset selection. This method keeps most of the information contained in individual assets while benefitting greatly from dimensionality reduction.

We base our model on a dynamic, characteristics-based factor model of returns. This kind of model was first studied by [Connor and Linton \(2007\)](#) and [Connor et al. \(2012\)](#), where they specified their model as:

$$y_{it} = \alpha_i + \sum_{j=1}^J g_j(X_{ji})f_{jt} + \varepsilon_{it}, \quad (3.1)$$

where  $y_{it}$  is the excess return on security  $i$  at time  $t$ ;  $f_{jt}$  is the  $j^{\text{th}}$  risk factor's return at time  $t$ ;  $X_{ji}$  is the  $j^{\text{th}}$  observable characteristic of firm  $i$ ;  $\alpha_i$  represents the intercept (mispricing) part of  $i^{\text{th}}$  asset return; and  $\varepsilon_{it}$  are the mean zero idiosyncratic shocks. They restricted characteristic-based loading  $g_j(\cdot)$  to be a uni-variate nonparametric function. To extend the dimension of the factor loading function  $g_j(\cdot)$ , Kelly et al. (2019) and Kim et al. (2019) specify both mispricing and factor loading parts as a parametric linear function of a large set of firm-specific characteristics as:

$$y_{it} = h(\mathbf{X}_i) + \sum_{j=1}^J g_j(\mathbf{X}_i) f_{jt} + \varepsilon_{it}. \quad (3.2)$$

They illustrated the validity of characteristics-based factor models and provided relevant empirical results. Ge et al. (2020) generalized the parametric part of Equation 3.2 as semiparametric functions to be consistent with earlier research. They also proposed power-enhanced tests to verify their model, concluding that the semiparametric mispricing component  $h(\mathbf{X}_i)$  was only significant during certain rolling windows.

Section two describes the econometric framework for our model. We assume the returns are generated by the asset pricing model in Ge et al. (2020) and that the factor risk premia are predictable based on a single-index function involving a set of both stationary and nonstationary predictors.

Section three presents the general portfolio management problem, and our restricted class of portfolio selection rules in which the problem is divided into two steps. In the first step, the investors choose a set of characteristics-based sub-portfolios that are well-diversified and mimic the returns of the underlying unobservable factors. In the second step, the investors choose a dynamic combination of these sub-portfolios and a risk-free asset dependent upon their time-varying information set and utility function. The information is specified as a single-index function, which is well-approximated by orthogonal series, allowing for both stationary and nonstationary covariates. We show that under reasonable conditions on risk preference, the two-step selection rule has asymptotically zero impact on an investor's expected utility as the number of assets grows to infinity, relative to the unattainable true optimal choice.

Section four derives estimators for both steps. In the factor-tilt step, the factor-mimicking portfolios are constructed by the linear combination of estimated characteristics-based factor loadings. To diversify the idiosyncratic shocks further, the weight for each factor loading function is estimated through a constrained quadratic objective function. In the second

step, called factor timing, the optimization of the expected utility function is solved using the continuously-updating GMM, as in [Hansen et al. \(1996\)](#). The weights allocated to the risk-free asset and sub-portfolios are determined by the single-index function approximated through Hermite polynomials, which allows for both stationary and nonstationary predictors, as in [Dong et al. \(2016b\)](#). The coefficients of those orthogonal bases are estimated by solving the sample counterpart under the continuously-updating GMM framework. Section five documents the hypothesis tests on the significance of these predictors included in the single-index function.

Section six presents the empirical findings. We apply our approaches to monthly CRSP and FRED data and reveal some popular predictive variables' nonstationarity and significance. Furthermore, we find our portfolios have different but outstanding performances under various levels of risk aversion. Finally, the results of the in-sample and out-sample are similar and reflect the risk preference of the investor.

Section seven concludes and discusses the paper, while proofs of theorems and supplementary tables are arranged in the Appendix.

## 3.2 The Model of Asset Returns

We assume that there is a large panel of monthly assets excess returns generated by the characteristics-based model:

$$y_{it} = \sum_{j=1}^J g_j(\mathbf{X}_i)(f_{jt} + \phi_{jt}) + \varepsilon_{it}, \quad (3.3)$$

where  $y_{it}$  is  $i^{th}$  stock's excess return at time  $t$  while  $\mathbf{X}_i$  is a large set of assets' P-vector of characteristics, which is regarded as time-invariant within a short time window;  $g_j(\mathbf{X}_i)$  is the  $j^{th}$  characteristics-based factor loading, which is specified as a multivariate additive semiparametric function, where  $g_j(\mathbf{X}_i) = \sum_{p=1}^P \mu_{jp}(X_{ijp})$  with  $\mu_{jp}(X_{ijp})$  the  $p^{th}$  univariate unknown characteristic function. The factor returns  $\mathbf{F}_t = (f_{1t}, \dots, f_{Jt})^\top$  are the common sources of risk in assets returns at time  $t$  with associated means  $\boldsymbol{\phi}_t = \{\phi_{1t}, \dots, \phi_{Jt}\}$ . The asset-specific return  $\varepsilon_{it}$  is conditional zero mean, i.e.,  $E(\varepsilon_{it}|\mathbf{X}_i) = 0$ .

This framework is an extension of [Connor and Linton \(2007\)](#) and [Connor et al. \(2012\)](#), who assumed the factor beta function  $g(\cdot)$  to be uni-variate. This model is a special case of [Ge et al. \(2020\)](#) by replacing the mispricing component with the mean value  $\phi_{jt}$  of the  $j^{th}$  risk factor.

We allow for time variation in the characteristics of assets across rolling windows. We treat the  $n \times P$  matrix of characteristics in the  $t^{th}$  rolling window  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^\top$ , as a random draw from a multivariate population distribution. Furthermore, the investor can observe  $\mathbf{X}$  before rolling block  $t$ , and then choose his time  $t$  portfolio.

We define the  $n \times J$  matrix  $G(\mathbf{X}) = (g_1(\mathbf{X}), \dots, g_J(\mathbf{X}))$  and  $g_j(\mathbf{X}) = (g_j(\mathbf{X}_1), \dots, g_j(\mathbf{X}_n))^\top$ , and the matrix form of the **demeaned** assets excess returns at time  $t$  is :

$$\mathbf{Y}_t = G(\mathbf{X})\mathbf{F}_t + \boldsymbol{\varepsilon}_t, \quad t = 1, 2, \dots, T, \quad (3.4)$$

where  $\mathbf{Y}_t$  is a  $n \times 1$  matrix of the demeaned assets excess returns at time  $t$ ,  $G(\mathbf{X}) = (g_1(\mathbf{X}), \dots, g_J(\mathbf{X}))$  is a  $n \times J$  factor loading matrix, and  $\boldsymbol{\varepsilon}_t$  is a  $n \times 1$  vector of asset-specific risks.

We allow for dynamic variation in the mean value of factor return premia. At the beginning of each period, a  $K \times 1$  vector of random signal  $\mathbf{z}_t = (z_{1t}, \dots, z_{Kt})^\top$  is observed by the investor before he chooses his portfolio. The expected return on the  $j^{th}$  factor in [Equation 3.3](#) is a nonlinear function of a fixed linear combination of these dynamic signals by coefficients  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^\top$  as:

$$\phi_{jt+1} = \pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t). \quad (3.5)$$

The vectors  $\mathbf{z}_t$  and  $\boldsymbol{\varepsilon}_t$  are assumed to be statistically independent. At each time  $t$ , the investor observes the characteristics of those assets  $\mathbf{X}$ , which is treated as time-invariant during this time block, and the dynamic signal  $\mathbf{z}_t$ . Then, the investor chooses his time  $t$  portfolio based on this information. Finally, at the  $(t + 1)^{th}$  period, his portfolio return depends upon the realized assets returns, which in turn depends on the realized factors returns  $\mathbf{F}_{t+1}$  and asset-specific returns  $\boldsymbol{\varepsilon}_{t+1}$ , respectively, according to [Equation 3.3](#).

### 3.3 A Two-Step Version of the Portfolio Choice Problem

This section first defines the utility function of a rational decision-maker and then describes how the optimal portfolio weights are chosen through a two-step procedure. In step one, the investor chooses characteristics-based factor-mimicking sub-portfolios based on a linear combination of the beta function  $\sum_{j=1}^J g_j(\mathbf{X})$  in [Equation 3.3](#). Step two combines these sub-portfolios optimally using expected utility as the investor's objective function, based on a dynamic index.



### 3.3.1 Utility Function of the Investor

The investor in our model is myopic. He chooses his portfolio for time  $t$  to maximize one-period expected utility of return. We assume that his return at time  $t$  is  $W_t$  and his risk-averse von Neumann Morgenstern preference is defined over  $W_t$  with a lower bound on the second derivative:

$$\frac{d}{dW}u(W) > 0, -c < \frac{d^2}{dW^2}u(W) < 0 \quad (3.6)$$

Additionally, we define the optimal portfolios weights  $n \times 1$  vector  $\mathbf{w}^*$  such that:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} E[u(r_{ft} + \mathbf{w}^{*\top} \mathbf{r}_t) | \mathbf{X}, \mathbf{z}_t], \quad (3.7)$$

where  $r_{ft}$  is the risk-free return and  $\mathbf{r}_t$  is a  $n \times 1$  vector of stock returns at time  $t$ . In practice, the optimal  $\mathbf{w}^*$  is hard to determine and unstable when  $n$  is large or the trading frequency is high, as discussed in the Introduction. Therefore, we consider optimal portfolio choice under a restriction on portfolio weights. Rather than choosing asset weights directly, the investor chooses a set of  $J$  characteristics-based portfolios to approximately mimic the factors. Then, in the second step, the investor combines these factor-mimicking sub-portfolios optimally using his expected utility function conditional on a group of predictors.

### 3.3.2 Step 1: Factor-mimicking Sub-portfolios

In this subsection, we propose a method to construct factor-mimicking sub-portfolios based on [Equation 3.3](#) and discuss the properties of these sub-portfolios.

We propose a semiparametric weighting function to mimic the risk factors  $\mathbf{F}_t$ , which is in the form of a linear combination of characteristics-based factor loadings as in [Equation 3.3](#):

$$b_j(\mathbf{X}_i) = \gamma_{j1}g_1(\mathbf{X}_i) + \cdots + \gamma_{jJ}g_J(\mathbf{X}_i), \quad (3.8)$$

therefore, the portfolio weight of  $i^{th}$  asset to construct the  $j^{th}$  sub-portfolio is  $\frac{1}{n}b_j(\mathbf{X}_i)$ .

The weighting matrix of assets to mimic all  $J$  factors is as follows:

$$B(\mathbf{X}_i) = \frac{1}{n} \mathbf{\Gamma} G(\mathbf{X}_i)^\top, \quad (3.9)$$

where

$$B(\mathbf{X}_i) = \frac{1}{n} \begin{pmatrix} b_1(\mathbf{X}_i) \\ b_2(\mathbf{X}_i) \\ \vdots \\ b_J(\mathbf{X}_i) \end{pmatrix},$$

$$\mathbf{\Gamma} = \begin{pmatrix} \gamma_{11} & \cdots & \gamma_{1J} \\ \gamma_{21} & \cdots & \gamma_{2J} \\ \cdots & \cdots & \cdots \\ \gamma_{J1} & \cdots & \gamma_{JJ} \end{pmatrix}, \quad G(\mathbf{X}_i)^\top = \begin{pmatrix} G_1(\mathbf{X}_i) \\ G_2(\mathbf{X}_i) \\ \vdots \\ G_J(\mathbf{X}_i) \end{pmatrix}.$$

Thus, the  $J \times 1$  factor-mimicking portfolio returns vector at time  $t$  is calculated as:

$$\mathbf{Q}_t(\mathbf{X}) = \begin{pmatrix} q_{1t}(\mathbf{X}) \\ q_{2t}(\mathbf{X}) \\ \vdots \\ q_{Jt}(\mathbf{X}) \end{pmatrix} = \sum_{i=1}^n B(\mathbf{X}_i) y_{it}. \quad (3.10)$$

The factor-mimicking portfolio vector has at least two attractive properties, which are listed as Theorems.

An investor who uses a semiparametric characteristics-based weighting function to choose sub-portfolios rather than individual assets  $i$  sacrifices the flexibility to weight assets differently based on the properties of their asset-specific returns  $\varepsilon_{it}$ , since the sub-portfolio weight function [Equation 3.10](#) only differentiates assets by their characteristic vectors. However, for both hedge fund managers and researchers, there are no satisfactory rules for choosing thousands of assets robustly. Furthermore, some weighting strategies have to be rebalanced once per trading day, and even more frequently for some strategies. This high-speed decision-making problem is intractable without some simplifying applicable rules like [Equation 3.10](#).

### 3.3.3 Step 2: Factor-timing Portfolio Based on Dynamic Signals

This subsection describes how to approximate the dynamic signal function  $\pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t)$  in [Equation 3.5](#), and how to use this function as dynamic weights assigned to those factor-mimicking sub-portfolios in [subsection 3.3.2](#), to reflect information about their over-performance/under-performance on a risk-adjusted basis. This subsection captures the particular "factor-timing" strategy used by the investor.

Here, we define the objective function as:

$$\arg \max_{\boldsymbol{\theta}} E[u(\alpha r_{ft+1} + \boldsymbol{\Pi}(\boldsymbol{\theta}^\top \mathbf{z}_t)^\top \mathbf{Q}_{t+1}(\mathbf{X}))], \quad (3.11)$$

subject to

$$\|\boldsymbol{\theta}\|_2 = 1 \text{ and } \theta_1 > 0$$

and

$$\alpha + \sum_{j=1}^J \pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t) = 1$$

where  $r_{ft}$  is the risk-free return at time  $t$  and  $\alpha$  is its portfolio weight, and

$$\boldsymbol{\Pi}(\boldsymbol{\theta}^\top \mathbf{z}_t) = (\pi_1(\boldsymbol{\theta}^\top \mathbf{z}_t), \dots, \pi_J(\boldsymbol{\theta}^\top \mathbf{z}_t))^\top.$$

The first restriction is for identification purposes while the second is for unit investment. We do not restrict short selling and leverage. Equation 3.11 is a transformation of the objective function Equation 3.7. The  $n \times 1$  vector of assets' returns  $\mathbf{r}_t$  in Equation 3.7 is replaced by the vector of sub-portfolios' returns  $\mathbf{Q}_{t+1}$  conditional on  $\mathbf{X}$ , which compacts the information of  $\mathbf{r}_t$  through observed characteristics  $\mathbf{X}$ . Similarly, the dynamic weights of each asset  $\mathbf{w}^*$  is substituted by the dynamic information function  $\boldsymbol{\Pi}(\boldsymbol{\theta}^\top \mathbf{z}_t)$ , which is the mean function for the risk factors  $\boldsymbol{\phi}_t$  as in Equation 3.3. In other words, the objective function Equation 3.7 is a transformation of the utility function Equation 3.11 by incorporating conditional variables  $\mathbf{z}_t, \mathbf{X}$ .

Our purpose is to maximize the conditional expectation of the investor's utility function. The investment allocation to the  $j^{\text{th}}$  factor-tilt sub-portfolio is determined by the  $j^{\text{th}}$  information indicator  $\pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t)$ , which is a single-index function, to avoid the problem of "curse of dimensionality" caused by fully nonparametric methods. We specify fixed linear combinations as information input in an unknown function  $\pi_j(\cdot)$ , as stated by Ait-sahalia and Brandt (2001), for at least two reasons. Statistically, this can achieve a better convergence rate for estimates, and economically, a uni-variate index value provides meaningful and convenient descriptions of current investment opportunities. Meanwhile, these index functions' effects on each sub-portfolio can be highly nonlinear, as documented by Ait-sahalia and Brandt (2001). Therefore, we do not specify the functional form of  $\pi_j(\cdot)$ , allowing a parametric index function to influence each sub-portfolio's weight nonparametrically.

## 3.4 Econometric Methods

This section gives Assumptions for estimation, shows the estimation of characteristics-based factor loading  $G(\mathbf{X})$ , illustrates procedures for estimating  $\mathbf{\Gamma}$  in Equation 3.9 and  $\mathbf{B}, \boldsymbol{\theta}$  in ??, and theoretical results.

### 3.4.1 Assumptions

**Assumption 3.4.1.** As  $n \rightarrow \infty$ , we have:

$$\frac{1}{n} \mathbf{Y}^\top \mathbf{Y} \rightarrow_P \mathbf{M}_Y,$$

$$\mathbf{F}^\top \mathbf{F} = \begin{bmatrix} d_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & d_J \end{bmatrix},$$

where  $\mathbf{M}_Y$  is a positive definite matrix;  $\mathbf{F} = (\mathbf{F}_1, \dots, \mathbf{F}_T)$  and  $d_j$  are distinct and positive entries.

**Assumption 3.4.2.** We define  $\lambda_{\min}(M)$  and  $\lambda_{\max}(M)$  as the largest and the smallest eigenvalues of matrix  $M$ , respectively. Additionally, we define  $C_{\min}$  and  $C_{\max}$  to be positive constants such that:

$$C_{\min} \leq \lambda_{\min}\left(\frac{1}{n} \boldsymbol{\Phi}^\top(\mathbf{X}) \boldsymbol{\Phi}(\mathbf{X})\right) < \lambda_{\max}\left(\frac{1}{n} \boldsymbol{\Phi}^\top(\mathbf{X}) \boldsymbol{\Phi}(\mathbf{X})\right) \leq C_{\max},$$

$$C_{\min} \leq \lambda_{\min}(E(\boldsymbol{\varepsilon}_t \boldsymbol{\varepsilon}_t^\top)) < \lambda_{\max}(E(\boldsymbol{\varepsilon}_t \boldsymbol{\varepsilon}_t^\top)) \leq C_{\max},$$

$$C_{\min} \leq \lambda_{\min}\left(\frac{1}{n} \mathbf{G}(\mathbf{X})^\top \mathbf{P} \mathbf{G}(\mathbf{X})\right) < \lambda_{\max}\left(\frac{1}{n} \mathbf{G}(\mathbf{X})^\top \mathbf{P} \mathbf{G}(\mathbf{X})\right) \leq C_{\max},$$

as  $n \rightarrow \infty$ , where  $\boldsymbol{\Phi}(\mathbf{X})$  is the matrix of B-spline basis of characteristics;  $\mathbf{P}$  is the  $n \times n$  projection smoother with  $\mathbf{P} = \boldsymbol{\Phi}(\mathbf{X})(\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}))^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top$ .

**Assumption 3.4.3.**

$$\frac{1}{n} \mathbf{G}(\mathbf{X})^\top \mathbf{G}(\mathbf{X}) \rightarrow_P \mathbf{I}_J,$$

as  $n \rightarrow \infty$ , and  $\mathbf{I}_J$  is a  $J \times J$  identity matrix.

**Assumption 3.4.4.** *We assume that there exists a neighbourhood of  $\boldsymbol{\theta}$ ,  $\mathbf{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma}_\theta) \subset \Theta$ , such that for any  $\boldsymbol{\theta}_0 \in \mathbf{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma}_\theta)$ ,  $\boldsymbol{\theta}^\top \mathbf{z}_t$  is always an  $I(1)$  process.*

Both Assumption 3.4.1 and 3.4.3 are similar to those in Fan et al. (2016), which are used to separately identify risk factors and factor loadings. Given the orthogonal bases of B-splines and uncorrelated or weakly correlated characteristics, Assumption 3.4.3 is mild. Assumption 3.4.2 shows the explanation power of characteristics on factor loadings is non-vanishing and implies that the asset specific risk caused by  $\boldsymbol{\varepsilon}_t$  is diversifiable. Assumption 3.4.4 precludes the co-integration of the signal function  $\boldsymbol{\theta}^\top \mathbf{z}_t$ , which fits the framework of Dong et al. (2016b).

### 3.4.2 Estimation of Characteristics-based Factor Loadings

The estimation of the characteristics-based loading  $G(\mathbf{X})$  is the same as in Ge et al. (2020), through the Projected-PCA proposed by Fan et al. (2016). The idea is to project the  $n \times T$  asset excess returns  $\mathbf{Y}$  onto the B-splines space spanned by  $\mathbf{X}$  through a  $n \times n$  projection matrix  $P$ , where  $P = \Phi(\mathbf{X})(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top$ , and then, we collect the projected returns  $\hat{\mathbf{Y}}$ . Furthermore, we perform PCA on  $\frac{1}{n} \hat{\mathbf{Y}} \hat{\mathbf{Y}}^\top$ . The  $\hat{G}(\mathbf{X})$  is attained as  $\sqrt{n}$  times the largest  $J$  eigenvectors of  $\frac{1}{n} \hat{\mathbf{Y}} \hat{\mathbf{Y}}^\top$ . Due to the property of the PCA, the Assumption 3.4.3 is satisfied.

### 3.4.3 Estimation of the First Step–Factor-tilt

We assume that the investor chooses  $\boldsymbol{\Gamma}$  based on the following objective function:

$$\hat{\boldsymbol{\Gamma}} = \arg \min_{\boldsymbol{\Gamma}} \sum_{j=1}^J E(b_j(\mathbf{X}_i)^2) \quad (3.12)$$

subject to

$$E[Q_t(\mathbf{X})Q_t(\mathbf{X})^\top] = \mathbf{I}_J,$$

where  $\mathbf{I}_J$  is a  $J \times J$  identity matrix.

In other words, we choose the linear combination coefficients  $J \times J$  matrix to maximize the spread of the portfolio weights, specifically by minimizing the expected sum of squared portfolio weights, in the class of semiparametric functions of the characteristics, subject to an orthogonality constraint on the vector of sub-portfolios' returns. These portfolios are an econometrically-derived variant of the widely popular Small-Minus-Big (SMB) and High-minus-Low (HML) portfolios designed by Fama and French (1993) to capture the size-related and value-related return factors. Fama and French (1993) did not minimize the

sum of squared portfolio weights as was done in Equation 3.12, but they instead set the portfolio weights using capitalization weights, which, in the highly diversified US equity market, have a very low sum-of-square relative to the number of assets. Fama and French (1993) did not explicitly impose the orthogonality condition applied in Equation 3.12, but, as they noted, they chose their size and value breakpoints so that the portfolio returns would have a very low correlation. The reason that we set the orthogonal constraint here is to diversify idiosyncratic risks further.

Next we show that Equation 3.12 can have a Lagrangian solution. After expanding the constraint and under the independence between  $\mathbf{F}_t$  and  $\boldsymbol{\epsilon}_t$ , we have:

$$\begin{aligned} E[\mathbf{Q}_t(\mathbf{X})\mathbf{Q}_t(\mathbf{X})^\top] &= E[(\frac{1}{n}\boldsymbol{\Gamma}G(\mathbf{X})^\top G(\mathbf{X})\mathbf{F}_t + \frac{1}{n}\boldsymbol{\Gamma}G(\mathbf{X})^\top \boldsymbol{\epsilon}_t)(\frac{1}{n}\mathbf{F}_t^\top G(\mathbf{X})^\top G(\mathbf{X})\boldsymbol{\Gamma}^\top + \frac{1}{n}\boldsymbol{\epsilon}_t^\top G(\mathbf{X})\boldsymbol{\Gamma}^\top)] \\ &= E(\boldsymbol{\Gamma}\frac{G(\mathbf{X})^\top G(\mathbf{X})}{n}\mathbf{F}_t\mathbf{F}_t^\top\frac{G(\mathbf{X})^\top G(\mathbf{X})}{n}\boldsymbol{\Gamma}^\top) + \frac{1}{n^2}E(\boldsymbol{\Gamma}G(\mathbf{X})^\top \boldsymbol{\epsilon}_t\boldsymbol{\epsilon}_t^\top G(\mathbf{X})\boldsymbol{\Gamma}^\top) \\ &\rightarrow \boldsymbol{\Gamma}\mathbf{M}^G E(\mathbf{F}_t\mathbf{F}_t^\top)\mathbf{M}^G\boldsymbol{\Gamma}^\top + \frac{1}{n^2}\boldsymbol{\Gamma}G(\mathbf{X})^\top E(\boldsymbol{\epsilon}_t\boldsymbol{\epsilon}_t^\top)G(\mathbf{X})\boldsymbol{\Gamma}^\top \end{aligned} ,$$

which is a quadratic form in  $\boldsymbol{\Gamma}$ .

As for the objective function, we have:

$$\sum_{j=1}^J E(b_j(\mathbf{X}_i)^2) = E(\boldsymbol{\Gamma}^2 G(\mathbf{X}_i)^2),$$

which is linear in  $\boldsymbol{\Gamma}^2$ .

Therefore, we write this constrained optimization problem of sample analogues in the Lagrangian form:

$$L(\boldsymbol{\Gamma}) = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\Gamma}G(\mathbf{X}_i)^\top G(\mathbf{X}_i)\boldsymbol{\Gamma}^\top - \boldsymbol{\Lambda}^\top \text{vec}((\frac{1}{T} \sum_{t=1}^T \mathbf{Q}_t(\mathbf{X})\mathbf{Q}_t(\mathbf{X})^\top) - \mathbf{I}_J), \quad (3.13)$$

where  $\boldsymbol{\Lambda}$  is the  $\frac{1}{2}J(J+1)$  vector of Lagrangian multipliers;  $G(\mathbf{X}_i) = (g_1(\mathbf{X}_i), \dots, g_J(\mathbf{X}_i))$  and  $\text{vec}$  is the vectarization of a matrix.

The optimal  $\boldsymbol{\Gamma}$  and associated Lagrangian multipliers will solve the first order conditions:

$$\begin{aligned} \frac{\partial L}{\partial \boldsymbol{\Gamma}} &= \mathbf{0}^{J \times J}, \\ \frac{\partial L}{\partial \boldsymbol{\Lambda}} &= \mathbf{0}^{\frac{1}{2}J(J+1)}. \end{aligned}$$

Meanwhile, we collect the estimate  $\hat{\Gamma}$  to obtain the factor-mimicking sub-portfolios' returns as:

$$\hat{\mathbf{Q}}_t(\mathbf{X}) = \begin{pmatrix} \hat{q}_{1t}(\mathbf{X}) \\ \hat{q}_{2t}(\mathbf{X}) \\ \vdots \\ \hat{q}_{Jt}(\mathbf{X}) \end{pmatrix} = \sum_{i=1}^n \hat{B}(\mathbf{X}_i) y_{it} = \frac{1}{n} \sum_{i=1}^n \hat{\Gamma} \hat{G}(\mathbf{X}_i)^\top y_{it}, \quad (3.14)$$

where  $\hat{G}(\mathbf{X}_i)$  is the consistent estimate of Equation 3.3 as in Ge et al. (2020), where they specify  $G(\mathbf{X}_i)$  as an additive semiparametric function of asset-specific characteristics.

### 3.4.4 Estimation of the Second Step–Factor-timing

The next step derives an estimator for the dynamic portfolio allocation weighting functions  $\Pi(\boldsymbol{\theta}^\top \mathbf{z}_{t-1})$ .

To facilitate our estimation procedures, we approximate those unknown functions  $\pi_j(\cdot)$  by orthonormal bases similar to Dong et al. (2016b). Their methods can allow the elements of the information vector  $\mathbf{z}_t$  to be nonstationary. As pointed by Gao and Phillips (2013a) and Gao and Phillips (2013b), conventional kernel estimation as in Brandt (1999) and Aït-sahalia and Brandt (2001) methods may not be workable due to the breakdown of the limit theory, when  $\mathbf{z}_t$  is a multivariate  $I(1)$  process. In practice, some time series predictors are likely to be nonstationary, like the unemployment rate, inflation and exchange rates, among other economic indicators. Therefore, we apply a similar method as in the Dong et al. (2016b) to validate a more comprehensive application of our model.

Suppose all the link functions  $\pi_j$  belong to  $L^2(\mathbb{R}) = \{f(x) : \int f^2(x) dx < \infty\}$ . The Hermite function sequence  $\{\mathcal{H}_i\}$  is an orthonormal basis in  $L^2(\mathbb{R})$ :

$$\mathcal{H}_i(x) = (\sqrt{\pi} 2^i i!)^{-1/2} H_i(x) \exp(-\frac{x^2}{2}), \quad i \geq 0, \quad (3.15)$$

where  $H_i(x)$  are Hermite polynomials orthogonal with density  $\exp(-x^2)$ . The orthogonality reads  $\int H_i(x) H_j(x) dx = \delta_{ij}$ , the Kronecker delta.

Therefore, any continuous function  $\pi_j(\cdot) \in L^2(\mathbb{R})$  can be expanded into a linear combination of orthogonal series:

$$\pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t) = \sum_{l=0}^{\infty} \beta_{jl} \mathcal{H}_l(\boldsymbol{\theta}^\top \mathbf{z}_t). \quad (3.16)$$

We keep the first  $L - 1$  terms and leave the rest as approximation residues:

$$\pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t) = \sum_{l=0}^{L-1} \beta_{jl} \mathcal{H}_l(\boldsymbol{\theta}^\top \mathbf{z}_t) + \psi(\boldsymbol{\theta}^\top \mathbf{z}_t), \quad (3.17)$$

where  $\psi(\boldsymbol{\theta}^\top \mathbf{z}_t)$  is the approximation residues.

Furthermore, all  $J$  dynamic indicator functions can be approximated (we set the same truncation parameters for all functions for the purposes of notation simplicity only):

$$\boldsymbol{\Pi}(\boldsymbol{\theta}^\top \mathbf{z}_t) = \mathbf{B} \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_t) + \boldsymbol{\Psi}(\boldsymbol{\theta}^\top \mathbf{z}_t), \quad (3.18)$$

where

$$\mathbf{B} = \begin{Bmatrix} \beta_{10} & \dots & \beta_{1(L-1)} \\ \dots & \dots & \dots \\ \beta_{J0} & \dots & \beta_{J(L-1)} \end{Bmatrix}, \quad \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_t) = \begin{Bmatrix} \mathcal{H}_0(\boldsymbol{\theta}^\top \mathbf{z}_t) \\ \dots \\ \mathcal{H}_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_t) \end{Bmatrix},$$

and  $\boldsymbol{\Psi}(\boldsymbol{\theta}^\top \mathbf{z}_t)$  is the approximation error.

Therefore, the objective function [Equation 3.11](#) is transformed through replacing  $\boldsymbol{\Pi}(\boldsymbol{\theta}^\top \mathbf{z}_t)$  by  $\mathbf{B} \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1})$  as:

$$(\hat{\boldsymbol{\alpha}}, \hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}) = \arg \max_{\boldsymbol{\alpha}, \mathbf{B}, \boldsymbol{\theta}} E[u(\alpha r_{ft} + (\mathbf{B} \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}))^\top \hat{\mathbf{Q}}_t(\mathbf{X}))], \quad (3.19)$$

subject to

$$\begin{aligned} \|\boldsymbol{\theta}\|_2 &= 1 \text{ and } \theta_1 > 0 \\ \alpha + \sum_{j=1}^J \boldsymbol{\beta}_j^\top \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) &= 1 \end{aligned}$$

where  $\hat{\mathbf{Q}}_t(\mathbf{X})$  is the estimate of sub-portfolios from [Equation 3.14](#). This is essentially the same semiparametric estimation problem analyzed by [Aït-sahalia and Brandt \(2001\)](#). The procedure relies on the profile estimation of the single-index function. We iterate the first order condition to convergence after choosing initial values arbitrarily.

The first order condition of the maximization with respect to  $\mathbf{B}, \boldsymbol{\theta}$  is:

$$E[\mathbf{M}_t] = E \left\{ \begin{array}{c} u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) \otimes \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \\ u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) (\mathbf{B} \mathcal{H}'_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \otimes \mathbf{z}_{t-1}) \\ u'(\cdot) r_{ft} \end{array} \right\} = \mathbf{0}_{J \times L + J \times K + 1},$$



where  $\mathcal{H}'_L$  and  $u'(\cdot)$  are the first derivatives of the truncated orthonormal series and the investor's utility function respectively.

There are  $J \times L + J \times K + 1$  moment conditions to maximize the objective function.

These moment conditions can be used to construct standard GMM problem as was done in Hansen (1982):

$$(\alpha, \hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}) = \arg \min_{\alpha, \mathbf{B}, \boldsymbol{\theta}} E[\mathbf{M}_t]^\top \mathbf{S} E[\mathbf{M}_t]$$

subject to

$$\begin{aligned} \|\boldsymbol{\theta}\|_2 &= 1 \text{ and } \theta_1 > 0 \\ \alpha + \sum_{j=1}^J \boldsymbol{\beta}_j^\top \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) &= 1 \end{aligned}$$

where  $\mathbf{S}$  is the optimal weighting positive definite matrix as  $\mathbf{S} = \text{cov}(\mathbf{M}_t)^{-1}$ .

Then, we substitute these moment conditions  $E[\mathbf{M}_t]$  with corresponding sample counterparts as:

$$\mathbf{m}_t = \left\{ \begin{array}{c} \frac{1}{T} \sum_{t=1}^T u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) \otimes \mathcal{H}_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \\ \frac{1}{T} \sum_{t=1}^T u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) (\mathbf{B} \mathcal{H}'_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \otimes \mathbf{z}_{t-1}) \\ u'(\cdot) r_{ft} \end{array} \right\} = \mathbf{0}_{J \times L + J \times K + 1}.$$

Similarly,

$$\hat{\mathbf{S}} = \left( \frac{1}{T} \sum_{t=1}^T \left\{ \begin{array}{c} u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) \otimes \mathcal{H}_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \\ u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) \mathbf{I}_{L \times L} (\mathbf{B} \mathcal{H}'_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \otimes \mathbf{z}_{t-1}) \\ u'(\cdot) r_{ft} \end{array} \right\} \left\{ \begin{array}{c} u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) \otimes \mathcal{H}_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \\ u'(\cdot) \hat{\mathbf{Q}}_t(\mathbf{X}) \mathbf{I}_{L \times L} (\mathbf{B} \mathcal{H}'_{L-1}(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) \otimes \mathbf{z}_{t-1}) \\ u'(\cdot) r_{ft} \end{array} \right\}^\top \right)^{-1}.$$

We substitute the sample analogues and  $\hat{\mathbf{S}}$  into the objective function, and estimate  $\hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}$ :

$$(\hat{\alpha}, \hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}) = \arg \min_{\alpha, \mathbf{B}, \boldsymbol{\theta}} \mathbf{m}_t^\top \hat{\mathbf{S}} \mathbf{m}_t, \quad (3.20)$$

subject to

$$\begin{aligned} \|\boldsymbol{\theta}\|_2 &= 1 \text{ and } \theta_1 > 0 \\ \alpha + \sum_{j=1}^J \boldsymbol{\beta}_j^\top \mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_{t-1}) &= 1 \end{aligned}$$

Furthermore, we substitute Equation 3.20 into the optimization iteration, which is called the continuously-updating estimator; details can be found in Dong et al. (2018).

### 3.4.5 Hypothesis Tests

This section introduces the hypothesis tests that help us to understand which index variables are important to guide the construction of factor-timing portfolios. We apply a Wald test to infer the significance of  $\theta_j$ . We have the null and alternative hypotheses as follows:

$$\mathcal{H}_0 : \mathbf{C}\boldsymbol{\theta} = \mathbf{0}_{D \times 1}, \quad \text{against} \quad \mathcal{H}_1 : \mathbf{C}\boldsymbol{\theta} \neq \mathbf{0}_{D \times 1},$$

where  $\mathbf{C}$  is a  $D \times K$  fix matrix indicating the number of constraints  $D$ .

We denote the value of the objective function Equation 3.20 under  $\hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}$  as  $\mathcal{V}(\hat{\mathbf{B}}, \hat{\boldsymbol{\theta}})$  while under the null hypothesis  $\mathcal{H}_0$  as  $\mathcal{V}(\hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}^*)$ .

Therefore, if the null hypothesis is correct, we have:

$$T(\mathcal{V}(\hat{\mathbf{B}}, \hat{\boldsymbol{\theta}}^*) - \mathcal{V}(\hat{\mathbf{B}}, \hat{\boldsymbol{\theta}})) \sim \chi^2(D), \quad (3.21)$$

where  $\chi^2(D)$  is the chi-square distribution with degree of freedom  $D$ . This method is a minimum- $\chi^2$  test, the purpose of which is to check the minimized values of the objective function Equation 3.20 after imposing some restrictions.

We reject the null hypothesis if the test statistic exceeds the critical value.

### 3.4.6 Theoretical Results

**Theorem 3.4.1.** *Each sub-portfolio in the  $J \times 1$  factor-mimicking vector defined by Equation 3.10 is a linear combination of risk factors  $f_{jt}$  directly.*

Theorem 3.4.1 implies that we can control the similarity between sub-portfolios and risk factors by adjusting coefficients matrix  $\Gamma$ , which provides us with considerable flexibility.

**Theorem 3.4.2.** *The returns of portfolio defined by Equation 3.10 have asymptotically zero idiosyncratic variance.*

Theorem 3.4.2 illustrates that portfolio returns of factor-mimick sub-portfolios can diversify the asset-specific returns completely as the number of assets goes to infinity.

**Theorem 3.4.3.** *The restricted optimal portfolio weighting function chosen by ?? gives an approximately optimal portfolio.*

Theorem 3.4.3 demonstrates that, as the number of assets  $n \rightarrow \infty$ , our two-step procedure is approximately equivalent to Equation 3.7, which is the completely unrestricted asset-by-asset portfolio optimization because these two methods give the same expected utility asymptotically.

## 3.5 Empirical Study

### 3.5.1 Data Description

#### Index Variables

We use the same index variable set as Ait-sahalia and Brandt (2001). These variables are all at a monthly frequency:

- The **Default Spread** is the yield difference between Moody's Baa and Aaa rated bonds, observed from 1967-07-01 to 2017-06-01 (600 months in total) denoted as DS.
- The **Term Spread** is the yield difference between 10 and 1 year government bonds, observed from 1967-07-01 to 2017-06-01 (600 months in total) denoted as TS.
- The **Trend** is the difference between the log of the current S&P 500 index level and the log of the average index level over the previous 12 months, observed from 1967-07-01 to 2017-06-01 (600 months in total).
- The **Dividend Yield**, also called Dividend-to-Price, is the sum of dividends paid on the S&P 500 index over the past 12 months divided by the current level of the index observed from 1967-07-01 to 2017-06-01 (600 months in total). We use the percentage natural logarithm form of Dividend Yield, denoted as Ln(DY%).
- The **Risk Free** rate is obtained from the Fama-French factor model's risk-free rate, observed from 1967-07-01 to 2017-06-01, denoted in the percentage form as RF%.

In Table 3.1, both "Trend" and "RF%" have small variation while "RF%" has some strong correlation with "TS" and "Ln(DY%)". Apart from these, we also find that all of the index variables are not symmetrically distributed, shown by the non-zero skewness. As for the kurtosis, the table indicates that outliers are quite common among these variables.

In [Table 3.2](#), we conclude the results of the unit root test and autocorrelation. After the Dickey-Fuller tests, we fail to reject the null hypotheses that there are no unit roots among all index variables, especially for the "Ln(DY%)". That can also be found from [Figure 3.4](#). In terms of autocorrelation, almost all of the index series present persistent autocorrelation even for lag nine, "Ln(DY%)" showing a strong signal of autocorrelation coefficient of 0.94, as shown in [Figure 3.4](#) and [Figure 3.5](#). However, "Trend" is an exception, where the autocorrelation decays to zero and is negative after lag ten as shown in [Figure 3.3](#). These test results verify the necessity of applying orthogonal series to approximate the single index function with nonstationary covariates, as in [subsection 3.3.3](#).

The data above was collected from the websites of FRED and Multpl.

### Monthly Stock Data

We collected monthly stock returns from CRSP and firms' characteristics from Compustat, from 1965 to 2017. We constructed 33 characteristics following the methods of [Freyberger et al. \(2020b\)](#). Details of these characteristics can be found in the appendix of [Ge et al. \(2020\)](#). We construct characteristics from fiscal year  $t - 1$  to explain stock returns between July of year  $t$  to June of year  $t + 1$ . Following [Hou et al. \(2015\)](#), we adjust returns of delisted stocks. The method that we apply to estimate the [Equation 3.4](#) is similar to [Ge et al. \(2020\)](#). We only include firms with at least three years of data in Compustat. The values of firm-specific characteristics are updated annually since most characteristic data are reported every year. We use rolling windows to accommodate time-varying characteristics-based loadings, and the risk factors are estimated correspondingly.

Our in-sample analysis's period is 50 years, from July 1967 to June 2017 (600 months).

### 3.5.2 In-sample Factor-mimicking Portfolios

This section presents portfolios that mimic the annually updated risk factors estimated through [Equation 3.3](#). In this study, we choose the number of unobservable factors in [Equation 3.3](#) to be three. In [Ge et al. \(2020\)](#), they compared the effects of the number of factors through a simulation study, concluding that underestimating the number of factors can be problematic. However, their discussions mainly focused on the estimation of mispricing functions. We only have four dynamic index variables, and therefore, we follow the renowned research of [Fama and French \(1993\)](#) to set the number of factors to be three. According to the literature, three factors can capture the essential common variation in asset excess returns.

The methods are introduced in [subsection 3.3.2](#), and we utilize all 600 months of data and construct three such portfolios every year, assuming the number of risk factors in [Equation 3.3](#) to be three. Then, we conclude the descriptive statistics of these three sub-portfolios in [Table 3.3](#). The constraints in estimation determine the zero mean and unit variance. As for the correlation between risk factors and those sub-portfolios using observations of all 600 months, our first step works very well because the diagonal elements of correlation are quite high while the off-diagonal elements are negligible. That demonstrates that each sub-portfolio imitates only the target risk factor's variation accurately and leaves the rest uncorrelated. The weights put on each asset for these sub-portfolios are calculated through a constrained optimization, which restricts the similarity between sub-portfolios and risk factors. Furthermore, during certain years, the sub-portfolios behaved in the opposite direction of the imitated factor, which can also influence the average correlation over 50 years. The annual correlation can be found in [Table B.2](#), where some negatively correlated periods are presented.

### 3.5.3 Utility Function

We utilize the classic Constant Relative Risk Aversion (CRRA) utility function to model function  $u(W)$  in [Equation 3.6](#):

$$u(W) = \begin{cases} \frac{W^{1-\xi}}{1-\xi} & \text{if } \xi > 1; \\ \ln(W) & \text{if } \xi = 1, \end{cases}$$

where  $\xi$  is an integer and  $\xi = W \frac{\partial^2 u(W)/\partial W^2}{\partial u(W)/\partial W}$ , measuring the level of risk aversion. Therefore, under this setting, the investor is risk-averse and tries to maximize his expected utility function through factor-mimicking and factor-timing portfolio strategy. The CRRA utility function is twice differentiable, which can further facilitate our optimization algorithm.

### 3.5.4 Selection of Truncation Number

The value of  $L$  in [Equation 3.17](#), which is the truncation number in polynomials, needs to be determined here. Unfortunately, to the best of our knowledge, there is no rule of thumb for the best choice of  $L$ . We refer to [Dong et al. \(2015\)](#) and [Dong et al. \(2016a\)](#), where the authors determined  $L$  according to the number of observations  $n$ . However, the  $n$  in this study ranged from 468 to 2928. After trading off the computational burden and approximation accuracy, we choose  $L$  to be four throughout the empirical study.

## 90A Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection

Table 3.1 Index Variable Summary

Index name	T	Descriptive Statistics							Correlation Matrix				
		Mean	Variance	Median	Max	Min	Skewness	Kurtosis	DS	TS	Trend	Ln(DP%)	RF%
DS	600	1.08	0.2	0.94	3.38	0.55	1.82	7.34	1.00				
TS	600	1.12	1.39	1.23	3.40	-3.07	-0.32	2.72	0.09	1.00			
Trend	600	0.03	0.01	0.05	0.22	-0.4	-1.24	5.63	-0.27	0.07	1.00		
Ln(DY%)	600	1.00	0.17	1.05	1.83	0.10	-0.09	2.07	0.46	-0.26	-0.12	1.00	
RF%	600	0.39	0.08	0.41	1.35	0.00	0.51	3.44	0.23	-0.68	-0.01	0.65	1.00

This table documents the descriptive statistics of the index variables that are used in this empirical study as well as the correlations among them. To be consistent with most of the literature, we use the percentage values of DP and RF.

Table 3.2 Tests Summary

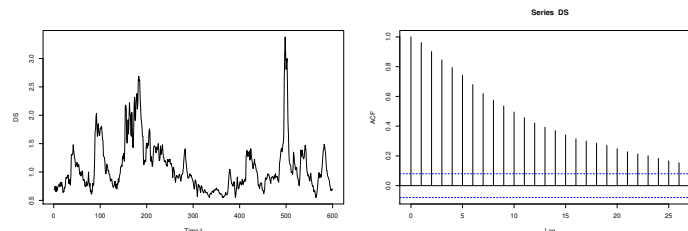
Index name	T	Unit root test				Autocorrelation		
		time trend	p-value (Trend)	$\Delta Y_t$	p-value ( $\Delta Y_t$ )	$\rho_3$	$\rho_6$	$\rho_9$
DS	600	$-2.3 \times 10^{-5}$ ( $2.8 \times 10^{-5}$ )	-0.81	$-3.82 \times 10^{-2}$ ( $1.11 \times 10^{-2}$ )	-3.43	0.85	0.68	0.54
TS	600	$9.27 \times 10^{-5}$ ( $7.51 \times 10^{-5}$ )	1.23	$-3.6 \times 10^{-2}$ ( $1.1 \times 10^{-2}$ )	-3.26	0.88	0.77	0.69
Trend	600	$4.1 \times 10^{-6}$ ( $8.75 \times 10^{-6}$ )	0.47	$-7.59 \times 10^{-2}$ ( $1.56 \times 10^{-2}$ )	-4.86	0.70	0.37	0.12
Ln(DY%)	600	$-1.8 \times 10^{-5}$ ( $1.25 \times 10^{-5}$ )	-1.44	$-8.86 \times 10^{-3}$ ( $5.16 \times 10^{-3}$ )	-1.72	0.98	0.96	0.94
RF%	600	$-6.6 \times 10^{-5}$ ( $2.05 \times 10^{-5}$ )	-3.22	$-5.48 \times 10^{-2}$ ( $5.16 \times 10^{-3}$ )	-4.23	0.94	0.90	0.87

This table summarizes the results of unit root tests and autocorrelations of those index variables. It reports the estimates, standard errors (in parentheses) and t-statistics of Dickey-Fuller test with trend individually. Autocorrelation column illustrates the correlation between the series and lag 3, lag 6 and lag 9 respectively, denoted as  $\rho_3$ ,  $\rho_6$ ,  $\rho_9$ .

Table 3.3 Factor-mimicking Portfolios Summary

Index name	T	Descriptive Statistics							Average Correlation		
		Mean	Variance	Median	Max	Min	Skewness	Kurtosis	$\hat{f}_1$	$\hat{f}_2$	$\hat{f}_3$
$\hat{q}_1$	600	0.00	1.00	0.08	2.81	-3.10	-0.38	3.33	0.48	-0.01	0.13
$\hat{q}_2$	600	0.00	1.00	0.03	2.70	-3.20	-0.15	3.42	-0.11	0.59	0.05
$\hat{q}_3$	600	0.00	1.00	0.05	2.60	-2.79	-0.10	2.82	0.06	-0.03	0.63

This table presents the descriptive statistics of factor-mimicking portfolios and their correlations with estimated risk factors.  $\hat{q}_1$ ,  $\hat{q}_2$ , and  $\hat{q}_3$  are constructed portfolios through all 600 months' data while  $f_1$ ,  $f_2$ , and  $f_3$  are three  $T \times 1$  factors estimated by rolling windows.



(a) Series Plot of DS      (b) Autocorrelation Plot of DS

Figure 3.1 The Plot of DS

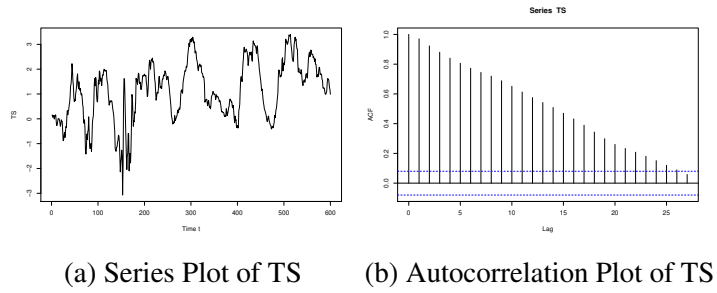


Figure 3.2 The Plot of TS

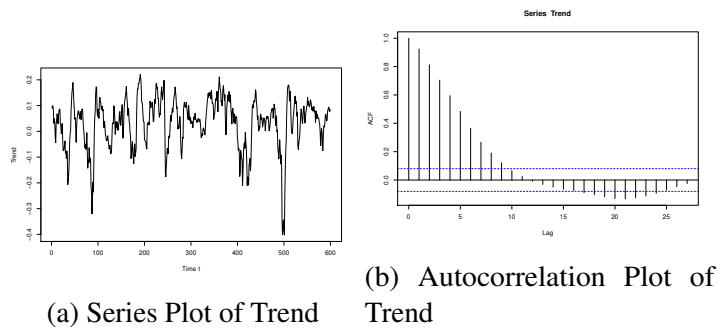


Figure 3.3 The Plot of Trend

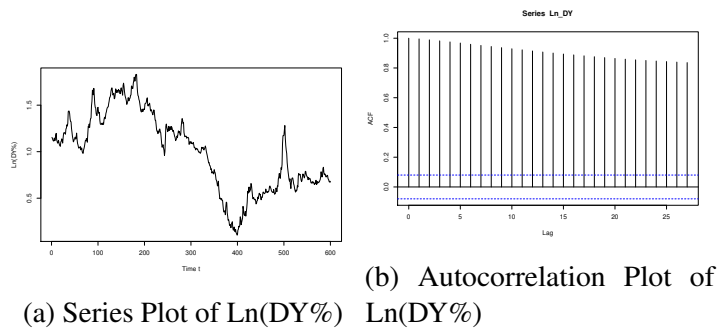


Figure 3.4 The Plot of Ln(DY%)

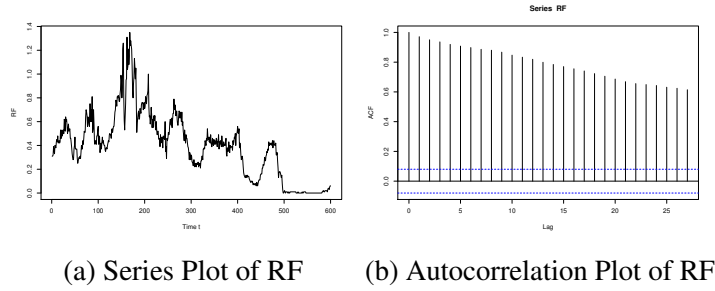


Figure 3.5 The Plot of RF

### 3.5.5 Estimation of Dynamic Signals

This section presents the estimation of single-index coefficient vector  $\theta$  and the results of the corresponding hypothesis tests. We use the CRRA utility function with various risk aversion levels  $\xi$ . Meanwhile, we also test the null hypothesis in subsection 3.4.5 to examine whether some of the coefficients of dynamic variables are significantly different from 0. These tests are used to show the importance of this dynamic information during the second step of portfolio management, namely factor-timing.

During our estimation, all the optimization processes converged, and the optimized values are reported. The in-sample results are based on the data for all 600 months and the estimation procedures are repeated under different risk-aversion levels  $\xi = 2, \xi = 5, \text{ and } \xi = 10$ . We then obtain the values of the objective function Equation 3.20, denoted as  $\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta})$ . The hypothesis tests are conducted by setting  $\theta_i = 0$ , where  $i$  indicates the  $i^{\text{th}}$  index variable. We denote the value of the objective function Equation 3.20 under  $\mathcal{H}_0 : \theta_i = 0$  as  $\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}_i = 0)$ . In addition,  $\chi^2$  statistics are calculated as  $T\Delta\mathbf{V} = T(\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}_i = 0) - \mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}))$ .

Table 3.4 Index Variable Summary

Index name	T	$\xi = 2$				$\xi = 5$				$\xi = 10$			
		$\hat{\theta}_i$	$\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}_i = 0)$	$\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta})$	$T\Delta\mathbf{V}$	$\hat{\theta}_i$	$\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}_i = 0)$	$\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta})$	$T\Delta\mathbf{V}$	$\hat{\theta}_i$	$\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}_i = 0)$	$\mathcal{V}(\hat{\mathbf{B}}, \hat{\theta})$	$T\Delta\mathbf{V}$
DS	600	0.15	0.23	0.006	134.4	0.06	0.0001	$2.5 \times 10^{-8}$	0.066	0.08	0.0026	$4.9 \times 10^{-10}$	1.56
TS	600	0.19	0.064	0.006	34.8	-0.34	0.0008	$2.5 \times 10^{-8}$	0.481	-0.06	$5.2 \times 10^{-9}$	$4.9 \times 10^{-10}$	0
Trend	600	0.03	0.01	0.006	2.4	0.06	$1 \times 10^{-6}$	$2.5 \times 10^{-8}$	0.006	0.03	$8.1 \times 10^{-10}$	$4.9 \times 10^{-10}$	0
Ln(DY%)	600	-0.97	0.06	0.006	32.4	-0.93	$3.2 \times 10^{-6}$	$2.5 \times 10^{-8}$	0.002	-0.995	$2 \times 10^{-8}$	$4.9 \times 10^{-10}$	0

This table reports the estimates and hypothesis test of dynamic index variables.  $\hat{\theta}_i$  is the estimate of the coefficient of the  $i^{\text{th}}$  index variable while  $\mathbf{V}$  represent the value of the objective function.  $\Delta\mathbf{V} = \mathcal{V}(\hat{\mathbf{B}}, \hat{\theta}_i = 0) - \mathcal{V}(\hat{\mathbf{B}}, \hat{\theta})$ .

The findings in Table 3.4 differ across the risk-aversion levels. When the magnitude of risk aversion is low, the influence of the dynamic index variables is significant. With  $\xi = 2$ , nearly all of the values of  $T\Delta\mathbf{V}$  exceed the 95% critical value of  $\chi^2(1)$ , which is 3.84, except for Trend. As the risk-aversion becomes larger, the importance of these dynamic variables



declines which can be confirmed when  $\xi = 5$  and  $\xi = 10$ , where all the four variables are insignificant. We compare the values of the objective function Equation 3.20, and most of them are quite similar and close to zero. That means the moment conditions in Equation 3.20 can be satisfied even if we restrict the coefficient of the  $i^{th}$  index variable to zero. Nevertheless, we cannot reject their joint significance.

### 3.5.6 In-sample Performance of Factor-timing Portfolios

This section presents portfolio performance estimated using in-sample data. As mentioned previously, there are two steps in constructing our dynamic portfolio, namely, factor-tilt and factor-timing steps. In subsection 3.5.2, we describe how to build the sub-portfolios that mimic the behavior of risk factors. This section solves the second step, factor-timing: choosing the time-varying weights for the risk-free asset and risky sub-portfolios. The dynamic weights are determined by a single-index function with a set of index variables. These variables capture investment opportunities. We standardize the amount of investment to be 1 unit and take the monthly returns as the wealth gleaned by the investor. We do not restrict leverage or short-selling to check the influence of the risk-aversion level  $\xi$ .

As we have 600 months in total, we record the average returns every year and annual standard deviations in Table 3.5 to save the space, and we calculate the Sharpe-ratio directly through  $mean(Return_t)/SD_{annual}$ . Table 3.5 shows the in-sample results from July 1967 to June 2017 under all three risk-aversion levels defined in subsection 3.5.3, and these results are compared with monthly S&P 500 returns.

Some findings here are substantial and worth discussing. Firstly, for investors who have relatively lower risk-aversion, the average portfolio returns are more rewarding, with some extremely high returns appearing as well. For example, when the risk-aversion level  $\xi = 2$ , the twelve-month average monthly returns can be 10.61 and 8.65. As for  $\xi = 10$ , the average monthly returns are more normal. Most monthly returns are around 5%, except for some outliers. Secondly, a higher risk-aversion level corresponds to more volatile returns, such as losing -2.33 monthly during the whole year when  $\xi = 2$ , provided the standard deviation of the monthly return is 6.44. But the circumstances can be much more favourable when  $\xi$  increases to 5 and 10, with the standard deviation of the monthly return being 3.98 and 2.81, respectively. Especially under  $\xi = 10$ , the returns are more stable. Thirdly, all of the portfolios under various  $\xi$  have a relatively low Sharpe-ratio, compared with S&P 500 returns, which may be due to the high volatility.

In this empirical study, we optimize over three risky sub-portfolios and one risk-free asset without restricting leveraging or short-selling, and the weight for each asset are plotted in [Figure 3.6](#). As we can see in (a), when the relative risk aversion level is low,  $\xi = 2$ , the weight for each sub-portfolio is variable, while the scale of the vertical axis here is broader than (b) and (c). As we increase  $\xi$ , the weights become more stable. Specifically, when the  $\xi$  increases to 10, the only substantial volatility in the weights appeared around the stock crash in March 2000.

### 3.5.7 Out-sample Performance of Factor-timing Portfolios

This section examines the out-sample performance of our two-step portfolio selection procedure. We test the last six months of the last ten years in our data set for various risk aversion levels. The coefficients of the dynamic information function are estimated using all of the past information while the sub-portfolios are estimated using the first six months each year. The "Return" in [Table 3.6](#) is calculated by substituting the predictors observed at the beginning of time  $t + 1$ . The sub-portfolios are constructed at time  $t$ , based on all the available data at the target year before time  $t + 1$ . [Table 3.6](#) also lists the assigned weights to each sub-portfolios and the risk-free asset using 1 unit of investment, represented by  $c_1, c_2, c_3$  and  $c_0$ . To summarize each column, we also provide the mean and standard deviation values at the end of the table, indicated by "ColMeans" and "ColStd".

As in [Table 3.6](#), most of the out-sample performance is quite similar to the in-sample performance in [Table 3.5](#). When the risk-aversion level is low such as  $\xi = 2$ , the variation of assets' weights is the largest and with an extensive range. Correspondingly, the realized monthly returns are also variable and high on average. The mean return of all 60 months is 0.36, which is very similar to that of the in-sample result 0.37. Not surprisingly, the out-sample standard deviation of 8.88 is bigger than that of the in-sample result (6.44).

When the risk-aversion level increases to  $\xi = 5$ , the weights' volatility decreases, and the mean return also falls from 0.36 to 0.27, which is similar to the in-sample result (0.29). Compared with the  $\xi = 2$  situation, the standard deviation of all the assigned weights and the monthly returns decline.

In the case of  $\xi = 10$ , all of the weights and monthly returns become more stable and less volatile. However, the average monthly return here is much lower than the in-sample result (0.17), with a smaller standard deviation of 1.51.

From the above analysis, we can conclude that our out-sample results are robust and vary according to the risk aversion levels. When the risk-aversion level is low, the investor

reassigns his weights broadly and frequently, with a high average monthly return but high volatility. As the risk aversion level increases, the investor adjusts his portfolios weights more moderately, and the monthly average return and its standard deviation are reduced.

Unfortunately, the extreme high volatility of all empirical results make a investor face the risk of bankruptcy every year, which means our method is practically inapplicable. We will discuss this problem in the next subsection.

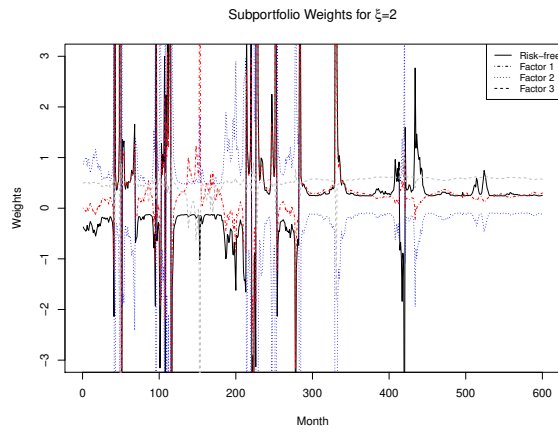
For this subsection, all the reported Sharpe-ratios are **not annualized**.

## 96A Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection

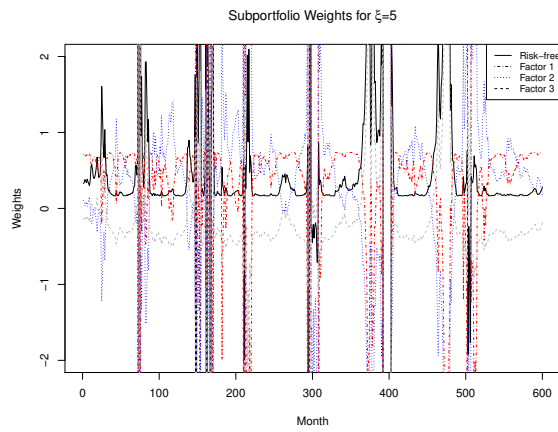
Table 3.5 Average Annual In-sample Results

n	$\xi = 2$			$\xi = 5$			$\xi = 10$			Average Monthly S&P 500		
	Return	SD	Sharpe-ratio	Return	SD	Sharpe-ratio	Return	SD	Sharpe-ratio	return	SD	Sharpe-ratio
468.00	-0.13	1.07	-0.12	0.15	0.70	0.22	0.06	0.61	0.10	0.01	0.03	0.27
894.00	-0.15	1.07	-0.14	0.20	0.47	0.41	0.08	0.63	0.13	-0.00	0.03	-0.03
1108.00	-0.13	1.11	-0.11	0.45	0.66	0.69	0.06	0.72	0.09	-0.02	0.04	-0.58
1199.00	-2.33	9.36	-0.25	0.10	0.88	0.12	0.04	0.59	0.08	0.02	0.03	0.87
1333.00	10.61	37.73	0.28	0.09	0.74	0.12	0.04	0.80	0.04	0.01	0.03	0.23
1409.00	-0.23	1.38	-0.16	0.10	0.77	0.13	0.06	0.82	0.08	-0.00	0.03	-0.08
1466.00	-0.16	1.03	-0.16	-1.02	3.50	-0.29	0.08	0.70	0.12	-0.01	0.04	-0.31
1560.00	-0.10	1.12	-0.09	0.35	1.19	0.30	0.05	0.70	0.07	0.00	0.07	0.07
1494.00	0.19	2.89	0.06	0.32	1.25	0.26	0.06	0.59	0.10	0.01	0.04	0.22
1292.00	0.52	7.48	0.07	0.14	1.29	0.11	0.03	0.60	0.05	-0.00	0.02	-0.09
1393.00	-0.12	0.73	-0.16	0.08	0.92	0.09	0.03	0.58	0.06	-0.00	0.03	-0.04
1340.00	-0.10	0.84	-0.12	0.38	0.94	0.41	0.05	0.79	0.06	0.00	0.03	0.12
1285.00	0.39	2.92	0.13	-0.13	5.27	-0.02	0.09	0.61	0.16	0.01	0.04	0.25
1181.00	-0.12	0.76	-0.16	-0.68	8.15	-0.08	0.03	0.57	0.06	0.01	0.03	0.47
1110.00	-0.15	0.83	-0.18	0.82	2.08	0.39	0.05	0.75	0.07	-0.01	0.04	-0.38
1044.00	-0.09	1.34	-0.07	-0.05	1.64	-0.03	-0.07	0.59	-0.11	0.04	0.03	1.05
1125.00	-0.73	1.72	-0.42	0.12	1.15	0.10	0.05	0.75	0.07	-0.01	0.02	-0.31
2192.00	-0.88	1.83	-0.48	-1.83	13.38	-0.14	0.05	0.73	0.06	0.02	0.03	0.59
2236.00	8.65	16.35	0.53	-0.87	2.35	-0.37	0.06	0.77	0.07	0.02	0.03	0.77
2273.00	0.07	0.53	0.13	0.08	0.87	0.09	0.05	0.95	0.05	0.02	0.03	0.57
2235.00	0.60	2.48	0.24	0.02	0.63	0.04	0.04	0.73	0.05	-0.01	0.06	-0.11
2270.00	-0.02	1.15	-0.02	0.28	0.95	0.30	0.09	0.58	0.16	0.02	0.02	0.68
2405.00	-0.33	1.19	-0.28	0.24	0.82	0.29	0.09	0.61	0.15	0.01	0.02	0.38
2376.00	2.03	3.87	0.53	0.02	1.00	0.02	0.03	0.63	0.05	0.01	0.05	0.11
2323.00	0.09	0.62	0.14	3.80	9.31	0.41	0.05	0.66	0.08	0.01	0.02	0.28
2344.00	0.05	0.68	0.08	2.45	3.21	0.76	0.03	0.61	0.05	0.01	0.01	0.56
2434.00	0.06	0.69	0.09	0.03	1.07	0.03	0.05	0.63	0.08	0.00	0.01	0.09
2548.00	-0.87	11.71	-0.07	-0.04	0.95	-0.04	0.10	0.56	0.17	0.01	0.02	0.79
2741.00	0.26	1.06	0.25	0.15	0.59	0.25	0.15	0.69	0.21	0.02	0.02	0.98
2928.00	0.10	0.55	0.19	0.14	0.74	0.18	0.25	0.53	0.47	0.02	0.04	0.64
2894.00	0.10	0.68	0.14	-0.03	1.90	-0.01	0.38	0.67	0.57	0.02	0.03	0.79
2905.00	0.11	0.73	0.16	-0.14	2.53	-0.05	0.25	0.64	0.39	0.02	0.05	0.33
2804.00	0.13	0.89	0.15	2.60	7.64	0.34	7.55	18.53	0.41	0.01	0.03	0.26
2570.00	0.10	0.69	0.14	-1.02	4.43	-0.23	0.95	1.64	0.58	-0.01	0.04	-0.34
2516.00	0.26	1.33	0.20	0.11	0.94	0.12	0.07	0.61	0.12	-0.02	0.05	-0.33
2491.00	0.11	0.89	0.13	0.07	1.05	0.07	0.03	0.60	0.05	-0.00	0.05	-0.01
2402.00	0.20	0.97	0.20	0.07	0.97	0.07	-0.04	0.58	-0.08	0.01	0.02	0.56
2326.00	0.06	0.90	0.07	0.01	0.50	0.03	0.11	0.58	0.18	0.01	0.02	0.23
2241.00	0.06	0.69	0.09	0.46	0.96	0.48	0.13	0.59	0.23	0.00	0.02	0.19
2178.00	0.12	0.73	0.17	7.04	15.93	0.44	0.15	0.60	0.25	0.02	0.02	0.89
2113.00	0.05	0.68	0.07	0.14	0.66	0.20	0.11	0.57	0.18	-0.01	0.04	-0.25
2023.00	-0.00	0.57	-0.00	-0.86	2.32	-0.37	0.04	0.85	0.04	-0.03	0.08	-0.33
2007.00	0.03	0.53	0.05	0.23	1.87	0.13	0.02	0.69	0.03	0.01	0.04	0.35
1924.00	-0.02	0.87	-0.02	0.02	1.14	0.02	-0.04	0.75	-0.05	0.01	0.02	0.62
1990.00	0.01	0.88	0.01	0.04	0.66	0.07	0.03	0.80	0.03	0.00	0.04	0.07
1937.00	0.00	0.89	0.00	-0.02	0.59	-0.03	0.01	0.79	0.01	0.02	0.02	0.77
1909.00	0.01	0.89	0.01	-0.02	0.94	-0.02	-0.01	0.70	-0.01	0.02	0.01	1.18
1872.00	0.00	0.69	0.01	-0.03	0.93	-0.04	-0.00	0.63	-0.00	0.01	0.02	0.31
1841.00	0.00	0.70	0.01	0.05	0.93	0.06	0.00	0.63	0.01	0.00	0.04	0.00
1826.00	0.01	0.70	0.02	-0.02	0.90	-0.02	-0.01	0.61	-0.01	0.01	0.01	0.94
Total mean	0.37	6.44	0.06	0.29	3.98	0.07	0.23	2.81	0.08	0.01	0.04	0.17

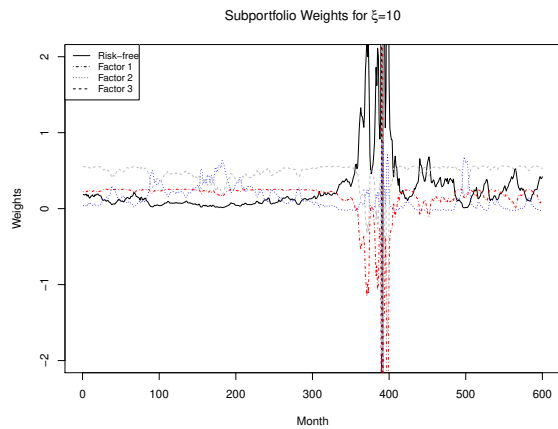
This table illustrates the in-sample results under various risk-aversion levels annually from July 1967- June 2017. n represents the number of stocks included in the portfolio. Both returns and standard deviations are calculated based on each year's results. The results of monthly S&P 500 returns are reported for comparison. No restrictions on leverage or short-selling.



(a) Monthly Weights Change Under  $\xi = 2$



(b) Monthly Weights Change Under  $\xi = 5$



(c) Monthly Weights Change Under  $\xi = 10$

Figure 3.6 The Plot of sub-portfolios Weights

# 98A Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection

Table 3.6 Monthly Out-sample Results Comparison

Year	$\xi = 2$					$\xi = 5$					$\xi = 10$				
	$c_0$	$c_1$	$c_2$	$c_3$	Return	$c_0$	$c_1$	$c_2$	$c_3$	Return	$c_0$	$c_1$	$c_2$	$c_3$	Return
2008	1.12	0.04	0.05	-0.21	-0.04	-0.14	0.59	0.22	0.32	0.12	0.07	0.73	0.01	0.20	0.65
	2.03	-0.13	-0.41	-0.49	0.21	-0.11	0.63	0.18	0.30	0.60	0.11	2.12	-0.37	-0.86	2.67
	2.41	-0.21	-0.64	-0.55	0.67	-0.23	0.89	-0.09	0.43	0.22	-0.14	-2.37	1.00	2.51	0.84
	2.31	-0.19	-0.58	-0.54	-0.82	-0.20	0.82	-0.02	0.40	0.69	-0.14	-2.31	0.99	2.46	4.20
	2.24	-0.17	-0.54	-0.53	-1.28	-0.14	0.69	0.12	0.33	1.35	0.62	11.15	-3.11	-7.65	-6.09
	1.85	-0.09	-0.31	-0.45	0.77	-0.11	0.63	0.18	0.30	0.49	0.11	2.02	-0.34	-0.79	3.93
2009	-0.00	1.47	0.79	-1.26	-1.60	0.15	-1.56	2.56	-0.15	6.52	-0.29	0.20	0.24	0.85	2.69
	-0.00	1.45	0.83	-1.28	0.73	0.10	-1.08	2.03	-0.06	-0.30	-0.26	0.31	0.03	0.93	0.31
	-0.00	1.46	0.81	-1.28	-1.44	0.11	-1.15	2.11	-0.07	0.92	-0.26	0.30	0.04	0.92	-0.44
	-0.00	1.49	0.77	-1.25	-0.86	0.06	-0.66	1.60	-0.00	-1.90	-0.25	0.39	-0.13	1.00	-1.87
	-0.00	1.45	0.83	-1.28	-0.05	0.03	-0.25	1.19	0.03	1.21	-0.26	0.67	-0.66	1.25	0.82
	-0.00	1.39	0.94	-1.33	-0.48	0.02	0.03	0.93	0.03	0.33	-0.39	1.86	-2.80	2.33	-0.63
2010	-0.00	0.22	0.02	0.75	-0.21	-0.01	2.06	-0.73	-0.31	-1.17	-0.00	2.17	-0.84	-0.33	-1.23
	-0.00	0.23	0.02	0.75	0.50	-0.01	2.04	-0.72	-0.31	-0.00	-0.00	2.18	-0.84	-0.34	-0.05
	-0.00	0.23	0.02	0.75	-0.72	-0.01	2.07	-0.74	-0.31	-2.10	-0.00	2.17	-0.84	-0.33	-2.17
	-0.00	0.22	0.03	0.75	0.31	-0.01	1.97	-0.67	-0.30	-2.02	-0.00	2.18	-0.84	-0.34	-2.31
	-0.00	0.27	-0.01	0.74	-0.86	-0.01	2.11	-0.78	-0.32	-1.49	-0.00	2.18	-0.83	-0.35	-1.51
	-0.00	0.30	-0.03	0.74	-0.71	-0.01	2.39	-1.03	-0.36	-1.00	-0.00	2.17	-0.84	-0.32	-0.97
2011	1.49	-0.90	0.61	-0.20	1.01	0.00	1.07	0.28	-0.35	-0.11	-0.02	0.37	0.66	0.00	0.36
	1.02	-0.91	1.10	-0.20	1.42	0.00	1.09	0.26	-0.35	-0.35	-0.02	0.37	0.64	0.01	0.12
	1.18	-0.92	0.96	-0.22	-1.19	0.00	1.06	0.28	-0.35	0.45	-0.03	0.39	0.62	0.02	0.09
	1.05	-0.92	1.07	-0.21	-1.56	0.00	1.08	0.27	-0.35	0.54	-0.03	0.39	0.62	0.02	0.05
	1.31	-0.92	0.83	-0.21	-2.11	0.00	1.04	0.31	-0.34	1.24	-0.03	0.40	0.61	0.02	0.73
	1.52	-0.89	0.56	-0.19	-0.78	0.00	1.01	0.33	-0.35	0.78	-0.03	0.43	0.56	0.04	0.70
2012	-33.41	5.01	8.21	21.20	53.27	0.04	0.75	0.17	0.04	1.64	0.19	0.12	0.30	0.39	0.96
	31.14	-4.95	-6.39	-18.80	-37.15	0.05	0.76	0.17	0.02	1.23	0.19	0.12	0.30	0.39	0.77
	22.59	-3.62	-4.46	-13.50	-6.39	0.05	0.76	0.17	0.03	0.30	0.17	0.12	0.33	0.39	0.16
	13.19	-2.17	-2.36	-7.67	13.97	0.05	0.77	0.17	0.01	-1.08	0.18	0.12	0.31	0.39	-0.93
	11.36	-1.88	-1.95	-6.53	9.04	0.07	0.79	0.17	-0.03	-0.80	0.23	0.12	0.26	0.38	-0.72
	14.33	-2.34	-2.61	-8.37	-4.71	0.07	0.79	0.17	-0.04	0.33	0.27	0.12	0.23	0.38	0.26
2013	1.06	-0.71	2.04	-1.40	0.44	0.05	1.18	-0.87	0.64	-0.30	0.09	0.19	0.36	0.36	-0.04
	1.40	-1.02	2.53	-1.91	-0.90	0.05	1.17	-0.85	0.63	0.43	0.09	0.19	0.36	0.36	0.03
	1.12	-0.76	2.12	-1.48	0.89	0.05	1.14	-0.77	0.58	-0.64	0.09	0.18	0.37	0.35	-0.17
	0.43	-0.12	1.10	-0.41	-0.26	0.05	1.12	-0.72	0.55	0.92	0.10	0.18	0.37	0.35	0.32
	1.09	-0.73	2.08	-1.43	-5.69	0.05	1.07	-0.60	0.48	2.90	0.10	0.17	0.40	0.33	0.15
	-0.42	0.65	-0.08	0.85	1.06	0.04	1.05	-0.53	0.44	0.94	0.10	0.17	0.40	0.33	0.46
2014	0.01	0.24	0.63	0.13	-1.20	-0.04	2.39	1.34	-2.68	-4.24	0.08	-0.76	1.23	0.45	-3.96
	0.01	0.25	0.63	0.12	0.02	-0.05	2.39	1.34	-2.68	0.29	0.09	-0.36	0.89	0.38	0.14
	0.01	0.25	0.63	0.12	0.35	-0.05	2.39	1.34	-2.68	0.43	0.09	-0.28	0.82	0.37	0.34
	0.01	0.25	0.63	0.12	0.75	-0.05	2.39	1.34	-2.68	1.40	0.09	-0.37	0.90	0.38	0.97
	0.01	0.26	0.62	0.11	1.29	-0.06	2.40	1.34	-2.69	1.29	0.11	0.00	0.59	0.29	1.15
	0.01	0.25	0.62	0.12	-0.01	-0.05	2.39	1.34	-2.68	0.33	0.10	-0.19	0.75	0.34	0.08
2015	0.37	0.36	1.09	-0.82	-1.66	-0.35	0.80	0.34	0.21	-0.13	0.40	0.20	0.29	0.12	-0.06
	0.50	0.54	1.42	-1.46	2.83	-1.35	4.17	-1.20	-0.62	4.29	0.37	0.20	0.24	0.20	-0.21
	0.51	0.55	1.43	-1.48	-0.58	-1.32	4.08	-1.16	-0.60	-0.60	0.36	0.20	0.22	0.22	-0.05
	0.42	0.43	1.21	-1.05	2.27	-0.56	1.51	0.01	0.04	3.29	0.39	0.19	0.27	0.15	1.19
	2.79	3.58	7.06	-12.44	-3.46	0.33	-1.46	1.36	0.77	-0.37	0.31	0.23	0.15	0.31	-0.09
	-1.06	-1.50	-2.36	5.92	2.20	0.17	-0.91	1.11	0.64	-0.17	0.28	0.27	0.10	0.35	-0.04
2016	1.26	-0.15	-0.86	0.75	-1.11	-0.01	1.45	-1.28	0.84	-0.05	0.00	0.09	-0.11	1.02	-0.86
	1.37	-0.29	-0.40	0.32	0.19	-0.01	1.12	-0.89	0.78	0.50	0.00	0.11	-0.10	0.99	0.37
	1.37	-0.30	-0.35	0.28	1.24	-0.01	1.09	-0.85	0.78	-1.09	0.00	0.11	-0.10	0.99	-0.01
	1.36	-0.28	-0.43	0.34	0.80	-0.01	1.03	-0.78	0.76	0.84	0.00	0.13	-0.10	0.97	1.02
	1.37	-0.29	-0.40	0.32	-0.18	-0.01	0.94	-0.67	0.74	-2.72	0.00	0.17	-0.09	0.92	-1.29
	1.37	-0.33	-0.19	0.15	-0.06	-0.01	0.85	-0.56	0.72	2.38	0.00	0.17	-0.09	0.92	0.18
2017	-0.00	-1.51	2.14	0.37	-0.92	0.60	-1.03	-2.37	3.80	-0.65	0.08	0.25	0.46	0.21	0.04
	-0.00	-1.43	2.06	0.37	0.23	0.59	-1.00	-2.32	3.73	0.01	0.10	0.21	0.51	0.18	0.23
	-0.00	-1.09	1.72	0.37	1.21	0.34	-0.24	-1.03	1.93	-0.14	0.14	0.10	0.65	0.11	0.53
	-0.00	-0.71	1.35	0.37	1.56	0.25	0.10	-0.40	1.06	0.52	0.21	-0.08	0.87	0.00	0.85
	-0.00	-0.60	1.23	0.37	1.18	0.24	0.16	-0.27	0.87	-0.08	0.26	-0.20	1.01	-0.07	1.18
	-0.00	-0.30	0.93	0.37	0.28	0.24	0.25	-0.02	0.54	0.07	0.60	-1.08	2.03	-0.55	0.45
ColMean	1.58	-0.18	0.47	-0.88	0.36	-0.02	0.95	0.03	0.04	0.27	0.08	0.53	0.15	0.25	0.07
ColStd	7.15	1.42	2.03	4.75	8.88	0.30	1.18	1.02	1.21	1.58	0.19	1.67	0.81	1.22	1.51

This table demonstrates the out-sample results under various risk-aversion levels of the last six months from 2008-2017.  $c_0$  represents the weight of the risk-free asset while  $c_1$ ,  $c_2$  and  $c_3$  show the weights of three sub-portfolios individually. "Return" represents the monthly return. "ColMean" and "ColStd" show the column means and standard deviations, respectively. No restrictions on leverage or short-selling.

### 3.5.8 Restrictive Sub-portfolios Weights

The volatile weights of sub-portfolios bring massive volatility to our dynamic investment strategy. In this subsection, we restrict the weights of risk-free asset and sub-portfolios to be within  $[0, 1]$  by a Cumulative Distribution Function of  $\mathcal{N}(0, \sigma^2)$ . This does not mean the restriction of leverage and short selling since the first step identifies optimal characteristics-based portfolios subject to orthogonal and unit variance constraints only without restrict each asset's weight. However, both steps restrict the magnitude of leverage and short-selling now.

For the weight of the  $j^{\text{th}}$  sub-portfolio  $\mathbf{q}_j$ , we have:

$$\pi_j(\boldsymbol{\theta}^\top \mathbf{z}_t) = \Phi(\mathbf{B}\mathcal{H}_L(\boldsymbol{\theta}^\top \mathbf{z}_t)), \quad j = 1, \dots, J,$$

where  $\Phi(\cdot)$  is the c.d.f of a normal distribution. The empirical results below set  $\sigma^2 = 20$ . The Sharpe-ratios increase enormously.

After controlling the leverage and short selling at the factor timing step, both returns and volatility drop dramatically, with a massive improvement of the annualized Sharpe-ratio. The attitudes to risks of investors also matter, which is reflected in [Figure 3.7](#). When the investor is less risk averse, nearly all his input are arranged to factor-mimicking portfolios, namely risky assets. However, when the risk aversion level increases to  $\xi = 10$ , the risk free asset takes up to 40% of the total investment.

To show the necessity of the first step, factor mimicking, we compare the performance of our portfolios in [Table 3.7](#) to the performance of portfolios using Fama French 3 (FF3) factors as factor-tilt sub-portfolios. In [Table 3.8](#), we summarise the results using the same factor timing step as in [Table 3.7](#), but in the first step, we use the FF3 factors as three sub-portfolios that are available to be invested by investors.

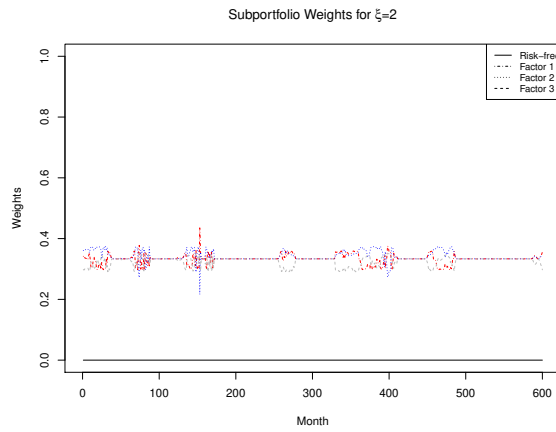
The FF3 portfolios enjoy slightly higher average annualized Sharpe-ratios, which are mainly due to the low volatility at all three risk aversion levels. However, the monthly return of FF3 portfolio is much lower, which can be also illustrated by [Figure 3.8](#), because the risk free asset takes a considerable amount for all levels of risk aversion.

Table 3.7 Average Annual In-sample Results

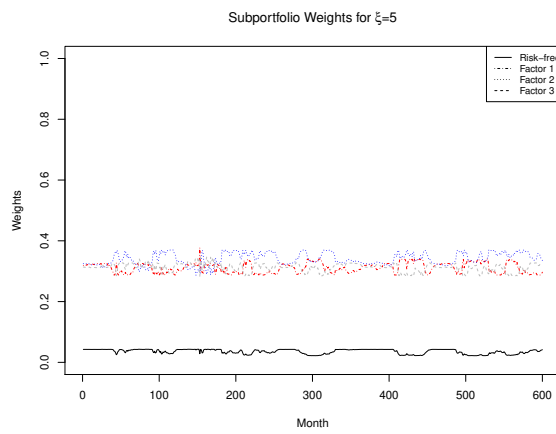
n	$\xi = 2$			$\xi = 5$			$\xi = 10$		
	Return	SD	Sharpe-ratio	Return	SD	Sharpe-ratio	Return	SD	Sharpe-ratio
468.00	0.40	0.53	2.66	0.39	0.50	2.69	0.24	0.31	2.70
894.00	-0.32	0.63	-1.73	-0.28	0.59	-1.67	-0.17	0.37	-1.63
1108.00	-0.42	0.72	-2.04	-0.40	0.68	-2.04	-0.25	0.43	-2.06
1199.00	0.28	0.54	1.81	0.25	0.52	1.67	0.21	0.36	2.01
1333.00	0.52	0.50	3.60	0.52	0.47	3.85	0.34	0.33	3.58
1409.00	-0.42	0.62	-2.35	-0.40	0.62	-2.26	-0.26	0.38	-2.36
1466.00	0.24	0.52	1.63	0.23	0.52	1.53	0.15	0.32	1.56
1560.00	0.22	0.58	1.32	0.22	0.54	1.41	0.17	0.40	1.47
1494.00	0.31	0.74	1.46	0.30	0.73	1.44	0.23	0.46	1.71
1292.00	-0.41	0.65	-2.18	-0.39	0.65	-2.06	-0.30	0.39	-2.63
1393.00	0.33	0.70	1.66	0.32	0.66	1.69	0.21	0.43	1.65
1340.00	-0.16	0.53	-1.04	-0.16	0.53	-1.07	-0.10	0.33	-1.05
1285.00	-0.10	0.75	-0.44	-0.08	0.74	-0.37	-0.04	0.47	-0.33
1181.00	0.38	0.54	2.47	0.38	0.51	2.60	0.25	0.33	2.66
1110.00	-0.69	0.31	-7.72	-0.68	0.29	-8.22	-0.43	0.18	-8.29
1044.00	-0.51	0.56	-3.17	-0.49	0.53	-3.21	-0.35	0.40	-3.04
1125.00	0.43	0.64	2.34	0.40	0.60	2.30	0.30	0.43	2.36
2192.00	0.10	0.31	1.11	0.15	0.29	1.77	-0.04	0.19	-0.80
2236.00	0.45	0.62	2.49	0.41	0.60	2.36	0.31	0.41	2.59
2273.00	-0.27	0.73	-1.28	-0.26	0.70	-1.28	-0.21	0.45	-1.58
2235.00	0.23	0.51	1.56	0.19	0.51	1.30	0.20	0.31	2.22
2270.00	0.45	0.62	2.49	0.43	0.59	2.50	0.27	0.38	2.49
2405.00	-0.37	0.69	-1.84	-0.37	0.64	-2.00	-0.22	0.41	-1.88
2376.00	0.04	0.73	0.17	0.04	0.70	0.19	0.05	0.49	0.35
2323.00	0.23	0.76	1.06	0.21	0.76	0.95	0.16	0.46	1.24
2344.00	0.57	0.42	4.70	0.57	0.42	4.67	0.41	0.30	4.71
2434.00	0.04	0.73	0.17	0.01	0.73	0.03	0.07	0.45	0.57
2548.00	0.63	0.65	3.37	0.60	0.60	3.48	0.40	0.38	3.59
2741.00	0.63	0.38	5.73	0.60	0.36	5.74	0.39	0.23	5.74
2928.00	0.03	0.69	0.14	-0.03	0.63	-0.15	-0.00	0.42	-0.01
2894.00	0.29	0.72	1.40	0.27	0.69	1.38	0.17	0.45	1.35
2905.00	0.26	0.57	1.60	0.26	0.54	1.67	0.17	0.34	1.69
2804.00	0.38	0.52	2.50	0.38	0.50	2.63	0.24	0.31	2.64
2570.00	0.38	0.66	1.98	0.36	0.65	1.93	0.23	0.41	1.97
2516.00	0.45	0.62	2.53	0.44	0.59	2.62	0.35	0.42	2.93
2491.00	0.54	0.45	4.15	0.51	0.45	3.94	0.40	0.29	4.91
2402.00	-0.61	0.40	-5.24	-0.59	0.40	-5.13	-0.43	0.27	-5.46
2326.00	0.49	0.58	2.94	0.48	0.56	2.93	0.33	0.38	2.94
2241.00	0.72	0.27	9.29	0.69	0.27	8.74	0.43	0.17	8.68
2178.00	0.72	0.52	4.74	0.68	0.51	4.67	0.43	0.32	4.67
2113.00	0.56	0.44	4.45	0.52	0.44	4.02	0.38	0.26	5.17
2023.00	0.25	0.55	1.61	0.22	0.55	1.36	0.23	0.33	2.40
2007.00	0.19	0.70	0.93	0.18	0.68	0.92	0.18	0.47	1.36
1924.00	0.59	0.47	4.35	0.58	0.49	4.11	0.39	0.25	5.26
1990.00	-0.60	0.47	-4.44	-0.59	0.45	-4.61	-0.43	0.31	-4.79
1937.00	0.63	0.38	5.74	0.61	0.36	5.83	0.42	0.25	5.85
1909.00	-0.34	0.69	-1.68	-0.31	0.68	-1.56	-0.26	0.47	-1.92
1872.00	-0.33	0.70	-1.61	-0.33	0.68	-1.69	-0.24	0.47	-1.79
1841.00	-0.34	0.56	-2.09	-0.37	0.55	-2.32	-0.25	0.37	-2.31
1826.00	0.21	0.74	0.98	0.19	0.71	0.93	0.14	0.47	1.01
Average	0.15	0.69	0.73	0.14	0.67	0.71	0.10	0.44	0.75

This table illustrates the in-sample results under various risk-aversion levels annually from July 1967- June 2017. n represents the number of stocks included in the portfolio. Both returns and standard deviations are calculated based on each year's results. **Sharpe-ratio is annualized. No leverage or short-selling for the factor timing step.**

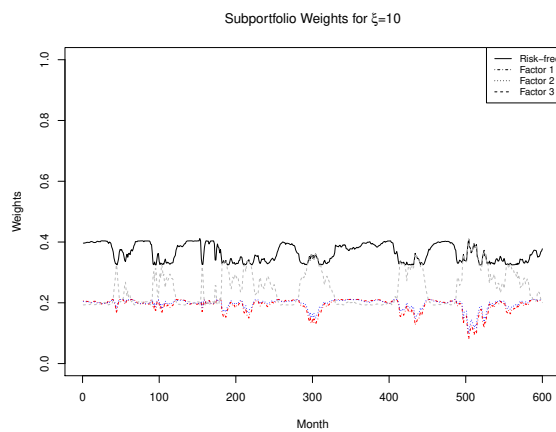




(a) Monthly Weights Change Under  $\xi = 2$



(b) Monthly Weights Change Under  $\xi = 5$



(c) Monthly Weights Change Under  $\xi = 10$

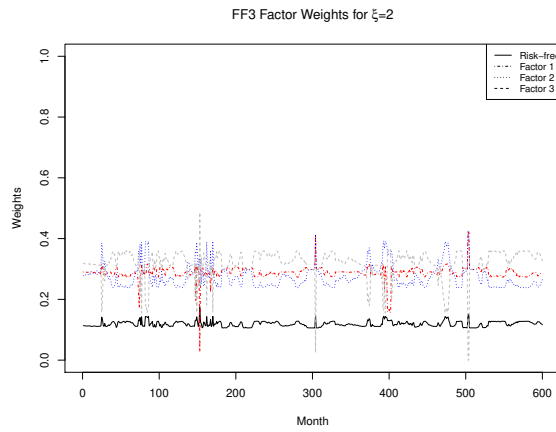
Figure 3.7 The Plot of Sub-portfolio Weights

## 102 Dynamic Semiparametric Characteristics-based Model for Optimal Portfolio Selection

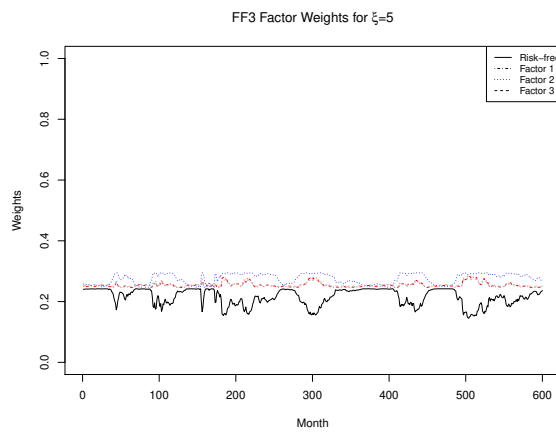
Table 3.8 Average Annual In-sample Results

n	$\xi = 2$			$\xi = 5$			$\xi = 10$		
	Return	SD	Sharpe-ratio	Return	SD	Sharpe-ratio	Return	SD	Sharpe-ratio
468.00	0.01	0.02	1.79	0.01	0.02	1.93	0.01	0.01	2.08
894.00	0.00	0.02	0.27	0.00	0.02	0.39	0.00	0.01	0.62
1108.00	-0.01	0.02	-1.78	-0.01	0.02	-1.77	-0.01	0.01	-1.44
1199.00	0.01	0.02	1.85	0.01	0.02	1.91	0.01	0.01	2.01
1333.00	-0.00	0.02	-0.14	0.00	0.02	0.09	0.00	0.01	0.13
1409.00	-0.00	0.01	-0.84	-0.00	0.01	-1.21	-0.00	0.01	-0.69
1466.00	-0.00	0.03	-0.04	0.00	0.02	0.22	0.00	0.02	0.44
1560.00	0.01	0.03	0.79	0.01	0.03	0.96	0.01	0.03	1.07
1494.00	0.01	0.03	1.20	0.01	0.03	1.18	0.01	0.02	1.34
1292.00	0.01	0.01	2.51	0.01	0.01	2.69	0.01	0.01	2.83
1393.00	0.01	0.01	1.99	0.01	0.01	2.65	0.01	0.01	2.61
1340.00	0.01	0.02	0.96	0.01	0.02	1.07	0.01	0.02	1.34
1285.00	0.01	0.02	1.12	0.00	0.02	0.54	0.00	0.02	0.90
1181.00	0.01	0.01	2.15	0.01	0.01	3.35	0.01	0.01	4.11
1110.00	-0.00	0.01	-0.27	0.00	0.01	0.37	0.00	0.01	1.50
1044.00	0.02	0.01	4.87	0.02	0.01	5.54	0.01	0.01	5.48
1125.00	0.00	0.01	0.90	0.00	0.01	0.77	0.00	0.01	1.49
2192.00	0.01	0.01	2.82	0.01	0.01	2.74	0.01	0.01	3.45
2236.00	0.01	0.01	2.20	0.01	0.01	2.59	0.01	0.01	2.79
2273.00	0.00	0.01	0.77	0.00	0.01	0.60	0.00	0.01	1.04
2235.00	0.00	0.02	0.31	0.00	0.02	0.30	0.00	0.02	0.44
2270.00	0.00	0.01	3.00	0.00	0.00	2.87	0.00	0.00	4.04
2405.00	-0.00	0.01	-1.31	-0.00	0.01	-0.92	-0.00	0.01	-0.48
2376.00	-0.00	0.02	-0.11	0.00	0.02	0.05	0.00	0.01	0.20
2323.00	0.01	0.02	1.47	0.01	0.01	1.53	0.01	0.01	1.63
2344.00	0.01	0.01	3.11	0.01	0.01	2.84	0.01	0.01	2.98
2434.00	0.00	0.01	0.95	0.00	0.01	0.99	0.00	0.01	1.15
2548.00	0.00	0.01	1.19	0.00	0.01	1.53	0.00	0.01	1.64
2741.00	0.00	0.01	1.81	0.01	0.01	2.01	0.00	0.01	2.42
2928.00	0.01	0.01	2.38	0.01	0.01	2.07	0.01	0.01	2.64
2894.00	0.01	0.01	1.73	0.01	0.01	2.01	0.01	0.01	2.30
2905.00	-0.00	0.02	-0.56	-0.00	0.02	-0.34	-0.00	0.02	-0.16
2804.00	0.00	0.03	0.14	0.00	0.01	0.12	0.00	0.01	0.43
2570.00	0.01	0.01	4.64	0.01	0.01	4.62	0.01	0.01	4.91
2516.00	0.00	0.02	0.68	0.00	0.02	0.85	0.00	0.01	0.92
2491.00	-0.00	0.02	-0.09	-0.00	0.02	-0.00	-0.00	0.02	-0.04
2402.00	0.01	0.01	2.56	0.01	0.01	2.60	0.01	0.01	2.58
2326.00	0.01	0.01	1.51	0.00	0.01	1.32	0.00	0.01	1.56
2241.00	0.01	0.01	1.54	0.00	0.01	1.67	0.00	0.01	1.89
2178.00	0.00	0.01	1.90	0.00	0.01	2.42	0.00	0.00	2.99
2113.00	-0.01	0.01	-2.35	-0.01	0.01	-2.03	-0.01	0.01	-1.94
2023.00	-0.00	0.04	-0.35	-0.01	0.04	-0.45	-0.01	0.04	-0.47
2007.00	0.01	0.03	1.07	0.01	0.03	1.11	0.01	0.02	1.09
1924.00	0.01	0.02	1.18	0.01	0.02	1.28	0.00	0.01	1.18
1990.00	-0.00	0.02	-0.18	-0.00	0.02	-0.18	-0.00	0.02	-0.20
1937.00	0.01	0.01	3.28	0.01	0.01	3.00	0.01	0.01	3.16
1909.00	0.01	0.01	1.59	0.01	0.01	1.51	0.00	0.01	1.53
1872.00	-0.00	0.01	-0.41	-0.00	0.01	-0.21	-0.00	0.01	-0.34
1841.00	-0.00	0.02	-0.36	-0.00	0.01	-0.47	-0.00	0.01	-0.42
1826.00	0.01	0.02	1.22	0.01	0.02	1.25	0.01	0.01	1.25
Average	0.00	0.02	0.73	0.00	0.02	0.81	0.00	0.01	0.97

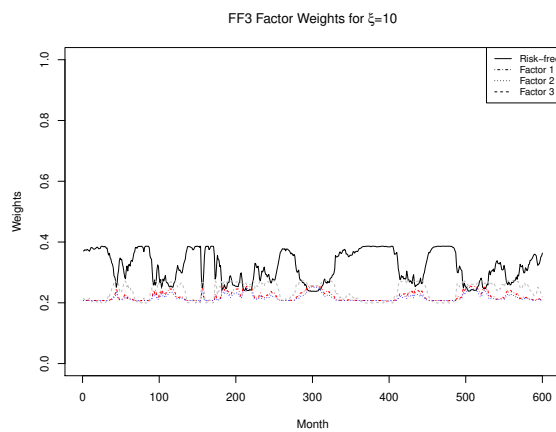
This table illustrates the in-sample results under various risk-aversion levels annually from July 1967- June 2017. n represents the number of stocks included in the portfolio. Both returns and standard deviations are calculated based on each year's results. **Sharpe-ratio is annualized. No leverage or short-selling for the factor timing step.**



(a) Monthly Weights Change Under  $\xi = 2$



(b) Monthly Weights Change Under  $\xi = 5$



(c) Monthly Weights Change Under  $\xi = 10$

Figure 3.8 The Plot of of FF 3 factors Weights

### **3.6 Conclusion**

This paper develops and tests a two-step portfolio selection procedure that relies on a large universe of investable assets and a set of dynamic predictors of factor-related returns. The first step in the procedure creates a collection of well-diversified mimicking portfolios to approximate the returns of pervasive risk factors. The second step uses a set of predictors including default spread, term spread, price trend, and dividend yield. These predictors are combined into a single-index function, which in turn determines a dynamic allocation of portfolio weights across the factor-mimicking portfolios in order to maximize investor's expected utility. Due to the nonstationarity of some predictive variables, we apply orthogonal series to approximate the single-index function in estimation. We apply the technique to fifty years of monthly U.S. data and find outstanding performance both in-sample and out-of-sample. We show empirically that the factor-mimicking portfolios have high correlation with the targeted factors and low correlation with others. Our dynamic portfolios perform well, both for high risk-aversion and low risk aversion investors, providing high average returns and also high return volatility for the less risk-averse and correspondingly lower average returns and lower volatility for the more risk-averse investor.

# Bibliography

- Y. Aït-sahalia and M. W. Brandt. Variable selection for portfolio choice. *The Journal of Finance*, 56(4):1297–1351, 2001.
- M. Ao, L. Yingying, and X. Zheng. Approaching mean-variance efficiency for large portfolios. *The Review of Financial Studies*, 32(7):2890–2919, 2019.
- M. W. Brandt. Estimating portfolio and consumption choice: A conditional euler equations approach. *The Journal of Finance*, 54(5):1609–1645, 1999.
- M. W. Brandt, P. Santa-Clara, and R. Valkanov. Parametric portfolio policies: Exploiting characteristics in the cross-section of equity returns. *The Review of Financial Studies*, 22(9):3411–3447, 2009.
- J. Y. Campbell and R. J. Shiller. Stock prices, earnings, and expected dividends. *The Journal of Finance*, 43(3):661–676, 1988.
- M. M. Carhart. On persistence in mutual fund performance. *The Journal of finance*, 52(1):57–82, 1997.
- J. Chen, D. Li, O. Linton, and Z. Lu. Semiparametric dynamic portfolio choice with multiple conditioning variables. *Journal of Econometrics*, 194(2):309–318, 2016.
- L. Chen, M. Pelger, and J. Zhu. Deep learning in asset pricing. *Available at SSRN 3350138*, 2020.
- X. Chen and D. Pouzo. Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals. *Econometrica*, 80(1):277–321, 2012.
- G. Connor and R. A. Korajczyk. Performance measurement with the arbitrage pricing theory: A new framework for analysis. *Journal of financial economics*, 15(3):373–394, 1986.
- G. Connor and O. Linton. Semiparametric estimation of a characteristic-based factor model of common stock returns. *Journal of Empirical Finance*, 14(5):694–717, 2007.
- G. Connor, M. Hagmann, and O. Linton. Efficient semiparametric estimation of the fama–french model and extensions. *Econometrica*, 80(2):713–754, 2012.
- D. R. Cox. Note on grouping. *Journal of the American Statistical Association*, 52(280):543–547, 1957.

- A. Deaton and J. Muellbauer. *Economics and consumer behavior*. Cambridge university press, 1980.
- C. Dong, J. Gao, and B. Peng. Semiparametric single-index panel data models with cross-sectional dependence. *Journal of Econometrics*, 188(1):301–312, 2015.
- C. Dong, J. Gao, and B. Peng. Another look at single-index models based on series estimation. *Available at SSRN 2858624*, 2016a.
- C. Dong, J. Gao, D. Tjøstheim, et al. Estimation for single-index and partially linear single-index integrated models. *The Annals of Statistics*, 44(1):425–453, 2016b.
- C. Dong, J. Gao, and O. B. Linton. High dimensional semiparametric moment restriction models. *Available at SSRN 3045063*, 2018.
- C. Dong, O. B. Linton, and B. Peng. A weighted sieve estimator for nonparametric time series models with nonstationary variables. *SSRN Working paper*, 2019.
- E. F. Fama and K. R. French. Business conditions and expected returns on stocks and bonds. *Journal of financial economics*, 25(1):23–49, 1989.
- E. F. Fama and K. R. French. Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics*, 33:3–56, 1993.
- E. F. Fama and K. R. French. A five-factor asset pricing model. *Journal of financial economics*, 116(1):1–22, 2015.
- J. Fan, Y. Liao, and M. Mincheva. Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society. Series B, Statistical methodology*, 75(4), 2013.
- J. Fan, Y. Liao, and J. Yao. Power enhancement in high-dimensional cross-sectional tests. *Econometrica*, 83(4):1497–1541, 2015.
- J. Fan, Y. Liao, and W. Wang. Projected principal component analysis in factor models. *Annals of statistics*, 44(1):219, 2016.
- Y. Fan and Q. Li. A kernel-based method for estimating additive partially linear models. *Statistica Sinica*, pages 739–762, 2003.
- G. Feng, S. Giglio, and D. Xiu. Taming the factor zoo. *Fama-Miller Working Paper*, 24070, 2017.
- W. D. Fisher. On grouping for maximum homogeneity. *Journal of the American statistical Association*, 53(284):789–798, 1958.
- J. Freyberger, A. Neuhierl, and M. Weber. Dissecting characteristics nonparametrically. Technical Report 5, 2020a.
- J. Freyberger, A. Neuhierl, and M. Weber. Dissecting characteristics nonparametrically. *The Review of Financial Studies*, 33(5):2326–2377, 2020b.

- J. Gao and C. Phillips. Functional coefficient nonstationary regression with non-and semi parametric cointegration. 2013a.
- J. Gao and P. C. Phillips. Semiparametric estimation in triangular system equations with nonstationarity. *Journal of Econometrics*, 176(1):59–79, 2013b.
- S. Ge, S. Li, and O. Linton. A dynamic network of arbitrage characteristics. *Available at SSRN 3638105*, 2020.
- S. Gu, B. Kelly, and D. Xiu. Empirical asset pricing via machine learning. *The Review of Financial Studies*, 33(5):2223–2273, 2020.
- L. P. Hansen. Large sample properties of generalized method of moments estimators. *Econometrica: Journal of the Econometric Society*, pages 1029–1054, 1982.
- L. P. Hansen, J. Heaton, and A. Yaron. Finite-sample properties of some alternative gmm estimators. *Journal of Business & Economic Statistics*, 14(3):262–280, 1996.
- T. Hastie and R. Tibshirani. Exploring the nature of covariate effects in the proportional hazards model. *Biometrics*, pages 1005–1016, 1990.
- E. Hjalmarrsson and P. Manchev. Characteristic-based mean-variance portfolio choice. *Journal of Banking & Finance*, 36(5):1392–1401, 2012.
- G. Hoberg and G. Phillips. Text-based network industries and endogenous product differentiation. *Journal of Political Economy*, 124(5):1423–1465, 2016.
- R. J. Hodrick. Dividend yields and expected stock returns: Alternative procedures for inference and measurement. *The Review of Financial Studies*, 5(3):357–386, 1992.
- K. Hou, C. Xue, and L. Zhang. Digesting anomalies: An investment approach. *The Review of Financial Studies*, 28(3):650–705, 2015.
- J. Huang, J. L. Horowitz, and F. Wei. Variable selection in non-parametric additive models. *Annals of Statistics*, 38(4):2282–2313, 2010a.
- J. Huang, J. L. Horowitz, and F. Wei. Variable selection in nonparametric additive models. *Annals of statistics*, 38(4):2282, 2010b.
- D. B. Keim and R. F. Stambaugh. Predicting returns in the stock and bond markets. *Journal of financial Economics*, 17(2):357–390, 1986.
- B. T. Kelly, S. Pruitt, and Y. Su. Instrumented principal component analysis. *Available at SSRN 2983919*, 2017.
- B. T. Kelly, S. Pruitt, and Y. Su. Characteristics are covariances: A unified model of risk and return. *Journal of Financial Economics*, 2019.
- S. Kim, R. A. Korajczyk, and A. Neuhierl. Arbitrage portfolios. *Georgia Tech Scheller College of Business Research Paper*, (18-43), 2019.
- A. B. Kock and D. Preinerstorfer. Power in high-dimensional testing problems. *Econometrica*, 87(3):1055–1069, 2019.

- O. Ledoit and M. Wolf. A well-conditioned estimator for large-dimensional covariance matrices. *Journal of multivariate analysis*, 88(2):365–411, 2004.
- O. Ledoit and M. Wolf. Nonlinear shrinkage of the covariance matrix for portfolio selection: Markowitz meets goldilocks. *The Review of Financial Studies*, 30(12):4349–4388, 2017.
- O. Ledoit, M. Wolf, et al. Nonlinear shrinkage estimation of large-dimensional covariance matrices. *The Annals of Statistics*, 40(2):1024–1060, 2012.
- C. Leng, Y. Lin, and G. Wahba. A note on the lasso and related procedures in model selection. *Statistica Sinica*, 16(4):1273, 2006.
- Q. Li. Efficient estimation of additive partially linear models. *International Economic Review*, 41(4):1073–1092, 2000.
- H. Liang, S. W. Thurston, D. Ruppert, T. Apanasovich, and R. Hauser. Additive partial linear models with measurement errors. *Biometrika*, 95(3):667–678, 2008.
- C. K. Liew. Inequality constrained least-squares estimation. *Journal of the American Statistical Association*, 71(355):746–751, 1976.
- Y. Lin and H. H. Zhang. Component selection and smoothing in smoothing spline analysis of variance models. *Annals of Statistics*, 34(5):2272–2297, 2006.
- O. Linton. Miscellaneous efficient estimation of additive nonparametric regression models. *Biometrika*, 84(2):469–473, 1997.
- O. Linton and W. Härdle. Estimation of additive regression models with known links. *Biometrika*, 83(3):529–540, 1996.
- O. B. Linton. Efficient estimation of generalized additive nonparametric regression models. *Econometric Theory*, 16(4):502–523, 2000.
- S. Ma and L. Yang. Spline-backfitted kernel smoothing of partially linear additive model. *Journal of Statistical Planning and Inference*, 141(1):204–219, 2011.
- H. M. Markowitz et al. Portfolio selection. *Journal of Finance*, 7(1):77–91, 1952.
- M. H. Pesaran and T. Yamagata. Testing capm with a large number of assets. In *AFA 2013 San Diego Meetings Paper*, 2012.
- D. Pollard. Strong consistency of k-means clustering. *The Annals of Statistics*, pages 135–140, 1981.
- D. Pollard et al. A central limit theorem for k-means clustering. *The Annals of Probability*, 10(4):919–926, 1982.
- S. Ross. The arbitrage theory of capital asset pricing. *Journal of Economic Theory*, 13(3):341–360, 1976.
- L. Schumaker. Spline functions: basic theory. *John Wiley&Sons, New York*, 1981.



- S. Sperlich, D. Tjøstheim, and L. Yang. Nonparametric estimation and testing of interaction in additive models. *Econometric Theory*, 18(2):197–251, 2002.
- W. Sun, J. Wang, Y. Fang, et al. Regularized k-means clustering of high-dimensional data and its asymptotic consistency. *Electronic Journal of Statistics*, 6:148–167, 2012.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- M. Vogt and O. Linton. Classification of non-parametric regression functions in longitudinal data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(1):5–27, 2017.
- H. Wang and C. Leng. A note on adaptive group lasso. *Computational Statistics & Data Analysis*, 52(12):5277–5286, 2008.
- L. Wang, L. Yang, et al. Spline-backfitted kernel smoothing of nonlinear additive autoregression model. *The Annals of Statistics*, 35(6):2474–2503, 2007.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- P. Zhao and B. Yu. On model selection consistency of lasso. *Journal of Machine learning research*, 7(Nov):2541–2563, 2006.
- H. Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429, 2006.



# Appendix A

## Proofs of Theorems

### A.1 Proofs

#### A.1.1 Proofs of Chapter 1

Let  $\boldsymbol{\beta}_{P_Z} = (\boldsymbol{\beta}_{P_1}^\top, \boldsymbol{\beta}_{P_2}^\top)$ ,  $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2)$ ,  $\beta_i$  is the  $i^{\text{th}}$  element of  $\boldsymbol{\beta}$ .  $\boldsymbol{\beta}_j$  is the  $j^{\text{th}}$  group of  $\boldsymbol{\beta}_{P_Z}$ , and  $\mathbf{X}_j$  is the covariates matrix of  $\mathbf{Z}$  in the second group.

In the first step, after applying KKT conditions, we obtain Lemma A.1.1 below.

**Lemma A.1.1.**

$$\frac{d\|\mathbf{Y} - \boldsymbol{\theta} - \mathbf{Z}\boldsymbol{\beta}\|^2}{d\beta_i} = \lambda_n \text{sign}(\hat{\beta}_i) \quad \text{for } \hat{\beta}_i \neq 0,$$

$$\frac{d\|\mathbf{Y} - \boldsymbol{\theta} - \mathbf{Z}\boldsymbol{\beta}\|^2}{d\beta_i} \leq \lambda_n \text{sign}(\hat{\beta}_i) \quad \text{for } \hat{\beta}_i = 0.$$

**Lemma A.1.2.** *Under Strong Irrepresentable Condition holds and a constant  $\eta > 0$ , then:*

$$P(\hat{\boldsymbol{\beta}}_{P_Z} =_s \boldsymbol{\beta}_{P_Z}) \geq P(E_A \cap E_B),$$

where:

$$E_A = \left\{ \frac{1}{\sqrt{n}} |(\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{z}_1^\top \mathbf{U}| < \sqrt{n} (|\boldsymbol{\beta}_{P_1}| - \frac{\lambda_n}{2n} |(\mathbf{V}_{Z_1 Z_1})^{-1} \text{sign}(\boldsymbol{\beta}_{P_1})|) \right\}$$

$$E_B = \left\{ \frac{1}{\sqrt{n}} |\mathbf{V}_{Z_2 Z_1} (\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top \mathbf{U} - \mathbf{Z}_2^\top \mathbf{U}| \leq \frac{\lambda_n}{2\sqrt{n}} \eta \right\},$$

The above equations hold for each entry.

The Lemma A.1.2 is borrowed from Proposition 1. of Zhao and Yu (2006). Proofs can be found in their Appendix.

**Proof of Theorem 1.4.1:** We give some notations before the proof. Let  $\boldsymbol{\tau} = \frac{1}{\sqrt{n}} (\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top \mathbf{U}$ , and  $\mathbf{v} = \frac{1}{\sqrt{n}} (\mathbf{V}_{Z_2 Z_1} (\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top \mathbf{U} - \mathbf{Z}_2^\top \mathbf{U})$ .

By Lemma A.1.2 we have:

$$1 - P(E_A \cap E_B) \leq P(E_A^c) + P(E_B^c) \leq \sum_{i=1}^{P_1} P(|\tau_i| \geq \sqrt{n} (|\beta_{P_1 i}| - \frac{\lambda_n}{2n} (\mathbf{V}_{Z_1 Z_1})^{-1} \text{sign}(\beta_{P_1 i}))) + \sum_{i=1}^{P_2} P(|v_i| \geq \frac{\lambda_n}{2\sqrt{n}} \eta_i).$$

Then we have

$$\mathbf{F}_\tau = \frac{1}{\sqrt{n}} (\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top,$$

therefore,

$$\mathbf{F}_\tau \mathbf{F}_\tau^\top = (\mathbf{V}_{Z_1 Z_1})^{-1}.$$

Given  $\lambda_{\min}(\mathbf{V}_{Z_1 Z_1}) > c_3$ , then we have  $\mathbf{V}_{Z_1 Z_1}^{-1} < c_5$  for each entry. Similarly, let

$$\mathbf{F}_v = \frac{1}{\sqrt{n}} (\mathbf{V}_{Z_2 Z_1} (\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top - \mathbf{Z}_2^\top),$$

and

$$\mathbf{F}_v \mathbf{F}_v^\top = \frac{1}{n} \mathbf{Z}_2^\top (\mathbf{I} - \mathbf{Z}_1^\top \mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top \mathbf{Z}_2.$$

Since  $\mathbf{I} - \mathbf{Z}_1^\top (\mathbf{V}_{Z_1 Z_1})^{-1} \mathbf{Z}_1^\top$  is idempotent, which only has the eigenvalues of 1 and 0, therefore  $\mathbf{F}_v \mathbf{F}_v^\top \leq c_4$  for each diagonal element.

Furthermore, we have:

$$\frac{\lambda_n}{n} |(\mathbf{V}_{Z_1 Z_1})^{-1} \text{sign}(\boldsymbol{\beta}_{P_1})| \leq \frac{c_5 \lambda_n}{n} \|\boldsymbol{\beta}_{P_1}\|_2$$

Given  $E(\varepsilon_i^{2k}) < \infty$ , then we have  $E(\tau_i^{2k}) < \infty$  and  $E(v_i^{2k}) < \infty$ . Therefore, the tail probability of  $\tau_i$  is bounded by:

$$P(\tau_i > T) = O(T^{-2k}),$$

furthermore, under  $\frac{\lambda_n}{\sqrt{n}} = o(n^{\frac{c_2-c_1}{2}})$ ,

$$\sum_{i=1}^{P_1} P(|\tau_i| \geq \sqrt{n}(|\beta_{P_1 i}| - \frac{\lambda_n}{2n}(\mathbf{V}_{Z_1 Z_1})^{-1} \text{sign}(\beta_{P_1 i}))) = P_1 O(n^{-kc_2}) = o(\frac{P_Z n^k}{\lambda_n^{2k}}). \quad (\text{A.1})$$

Similarly,

$$\sum_{i=1}^{P_2} P(|v_i| \geq \frac{\lambda_n}{2\sqrt{n}} \eta_i) = P_2 O(\frac{n^k}{\lambda_n^{2k}}) = o(P_Z \frac{n^k}{\lambda_n^{2k}}). \quad (\text{A.2})$$

Then, combining Equation A.1 and Equation A.2 gives Theorem 1.4.1.  $\square$

After grouping all the coefficients from step 1, we use  $\beta_j$  to represent the  $j^{\text{th}}$  group of  $\beta_{P_Z}$ .

we apply the KKT conditions again to obtain the Lemma A.1.3

**Lemma A.1.3.**

$$\frac{d\|\mathbf{Y} - \boldsymbol{\theta} - \mathbf{Z}\boldsymbol{\beta}\|^2}{d\boldsymbol{\beta}_j} = \hat{\omega}_j \tilde{\lambda}_n \frac{\hat{\boldsymbol{\beta}}_j}{\|\hat{\boldsymbol{\beta}}_j\|_2} \quad \text{for } \|\hat{\boldsymbol{\beta}}_j\|_2 \neq 0,$$

$$\frac{d\|\mathbf{Y} - \boldsymbol{\theta} - \mathbf{Z}\boldsymbol{\beta}\|^2}{d\boldsymbol{\beta}_j} \leq \hat{\omega}_j \tilde{\lambda}_n \quad \text{for } \|\hat{\boldsymbol{\beta}}_j\|_2 = 0,$$

Similar to Lemma 5 and Lemma 6 of Huang et al. (2010a), we give the following Lemmas:

**Lemma A.1.4.** Under Assumptions 1.4.1-1.4.4 and Condition 1.4.1-1.4.2:

$$P(\|\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j\|_2 \geq \|\boldsymbol{\beta}_j\|_2, \exists \mathbf{X}_j \in \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}) \rightarrow 0.$$

**Lemma A.1.5.** Under Assumptions 1.4.1-1.4.4 and Condition 1.4.1-1.4.2:

$$P(\|\mathbf{X}_j^\top (\mathbf{Y} - \mathbf{Z}_1 \boldsymbol{\beta}_1)\|_2 > \tilde{\lambda}_n \hat{\omega}_j / 2, \exists \mathbf{X}_j \notin \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}) \rightarrow 0$$

Proofs of Lemma A.1.4 and Lemma A.1.5 can be found in the Appendix of Huang et al. (2010a).

**Proof of Theorem 1.4.2 :** Theorem 1.4.2 satisfies the Condition 1 of Huang et al. (2010a).

Under Theorem 3.4.1, and Lemma A.1.3, we set  $\boldsymbol{\zeta} = (\frac{\hat{\omega}_j \hat{\boldsymbol{\beta}}_j}{2\|\hat{\boldsymbol{\beta}}_j\|})$ , for  $\mathbf{X}_j \in \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}$ .

Therefore, we have:

$$\hat{\boldsymbol{\beta}}_{P_1} = (\mathbf{Z}_1^\top \mathbf{Z}_1)^{-1} \mathbf{Z}_1^\top (\mathbf{Y} - \tilde{\lambda}_n \boldsymbol{\zeta}).$$

To proof Theorem 1.4.2, equivalently, we need to proof:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{P_1} &= \boldsymbol{\beta}_{P_1} \\ \|\mathbf{Z}_j^\top (\mathbf{Y} - \mathbf{Z}_1 \boldsymbol{\beta}_{P_1})\|_2 &\leq \tilde{\lambda}_n \hat{\omega}_j / 2 \quad \forall j \notin \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}. \end{aligned}$$

This is equivalently to show:

$$\begin{aligned} \|\boldsymbol{\beta}_j\|_2 - \|\hat{\boldsymbol{\beta}}_j\|_2 &< \|\boldsymbol{\beta}_j\|_2 \quad \forall j \in \mathcal{L} \cup \mathcal{R} \cup \mathcal{S} \\ \|\mathbf{Z}_j^\top (\mathbf{Y} - \mathbf{Z}_1 \boldsymbol{\beta}_{P_1})\|_2 &\leq \tilde{\lambda}_n \hat{\omega}_j / 2 \quad \forall j \notin \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}. \end{aligned}$$

Therefore,

$$\begin{aligned} P(\hat{\boldsymbol{\beta}}_{P_2} \neq \boldsymbol{\beta}_{P_2}) &\leq P(\|\boldsymbol{\beta}_j\|_2 - \|\hat{\boldsymbol{\beta}}_j\|_2 \geq \|\boldsymbol{\beta}_j\|_2 \quad \exists j \in \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}) \\ &\quad + P(\|\mathbf{Z}_j^\top (\mathbf{Y} - \mathbf{Z}_1 \boldsymbol{\beta}_{P_1})\|_2 > \tilde{\lambda}_n \hat{\omega}_j / 2 \quad \exists j \notin \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}). \end{aligned}$$

Theorem 1.4.1 shows

$$\hat{\omega}_j \rightarrow \infty, \quad \forall j \notin \mathcal{L} \cup \mathcal{R} \cup \mathcal{S},$$

where  $\omega_j$  is the specific penalty parameter of the  $j^{\text{th}}$  coefficient group.

Then,

$$P(\|\mathbf{Z}_j^\top (\mathbf{Y} - \mathbf{Z}_1 \boldsymbol{\beta}_{P_1})\|_2 > \tilde{\lambda}_n \hat{\omega}_j / 2, \quad \exists j \notin \mathcal{L} \cup \mathcal{R} \cup \mathcal{S}) \rightarrow 0$$

Therefore, under Lemma A.1.4 and Lemma A.1.5, the Theorem 1.4.2 follows.  $\square$

## A.1.2 Proofs of Chapter 2

Throughout the proofs, we have the number of observations  $n \rightarrow \infty$ , and time  $T$  is fixed.

**Proof of Theorem 2.6.1 :** In equation 5, we have

$$\mathbf{Y} = (\boldsymbol{\Phi}(\mathbf{X})\mathbf{A} + \boldsymbol{\Gamma} + \mathbf{R}^\mu(\mathbf{X}))\mathbf{1}_T^\top + (\boldsymbol{\Phi}(\mathbf{X})\mathbf{B} + \boldsymbol{\Lambda} + \mathbf{R}^\theta(\mathbf{X}))\mathbf{F}^\top + \mathbf{U},$$

Multiply time-demeaned matrix  $\mathbf{D}_T$  on both sides, where  $\mathbf{D}_T = \mathbf{I}_T - \frac{1}{T} \mathbf{1}_T \mathbf{1}_T^\top$ . Given time-invariant mispricing components, we obtain:

$$\mathbf{YD}_T = (\Phi(\mathbf{X})\mathbf{B} + \Lambda + \mathbf{R}^\theta(\mathbf{X}))\mathbf{F}^\top \mathbf{D}_T + \mathbf{UD}_T.$$

On-wards, we define  $\mathbf{YD}_T = \tilde{\mathbf{Y}}$  and  $\mathbf{F}^\top = \mathbf{F}^\top \mathbf{D}_T$ . Time-demeaned factors do not change their properties.

Next, multiple both sides by  $\mathbf{P} = \Phi(\mathbf{X})(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top$ ,

$$\hat{\mathbf{Y}} = (\Phi(\mathbf{X})\mathbf{B} + \mathbf{P}\Lambda + \mathbf{P}\mathbf{R}^\theta(\mathbf{X}))\mathbf{F}^\top + \mathbf{PUD}_T.$$

We decompose:

$$\mathbf{P}\tilde{\mathbf{Y}} = \hat{\mathbf{Y}} = \Phi(\mathbf{X})\mathbf{B}\mathbf{F}^\top + \mathbf{P}\Lambda\mathbf{F}^\top + \mathbf{PUD}_T + \mathbf{P}\mathbf{R}^\theta(\mathbf{X})\mathbf{F}^\top = \mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4,$$

as  $n \rightarrow \infty$  and  $n^v \rightarrow \infty$ , approximation error  $\mathbf{R}^\theta(\mathbf{X}) \rightarrow_P \mathbf{0}$  as in [Huang et al. \(2010b\)](#). Thus,  $\mathbf{e}_4^\top \rightarrow^P \mathbf{0}$ .

Under Assumption 2.2.1, we have following results:

for  $\frac{1}{n} \sum_{j=1}^3 \mathbf{e}_2^\top \mathbf{e}_j$ ,

$$\frac{1}{n} \mathbf{P}\Lambda \rightarrow^P \mathbf{0},$$

therefore,

$$\frac{1}{n} \sum_{j=1}^3 \mathbf{e}_2^\top \mathbf{e}_j + \frac{1}{n} \sum_{j=1}^3 \mathbf{e}_j^\top \mathbf{e}_2 \rightarrow^P \mathbf{0}.$$

For  $\frac{1}{n} \sum_{j=1}^3 \mathbf{e}_3^\top \mathbf{e}_j$ ,

$$\frac{1}{n} \mathbf{P}\mathbf{U} \rightarrow^P \mathbf{0},$$

therefore,

$$\frac{1}{n} \sum_{j=1}^3 \mathbf{e}_2^\top \mathbf{e}_j + \frac{1}{n} \sum_{j=1}^3 \mathbf{e}_j^\top \mathbf{e}_2 \rightarrow^P \mathbf{0}.$$

And only  $\frac{1}{n} \mathbf{e}_1^\top \mathbf{e}_1$  left, namely,

$$\frac{1}{n} \mathbf{e}_1^\top \mathbf{e}_1 = \mathbf{F} \frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n} \mathbf{F}^\top.$$

Under Assumption 2.6.1-2.6.3 and fixed  $T$ . A much smaller  $T \times T$  matrix  $\frac{1}{n} \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$  can be solved by asymptotic principal component by Connor and Korajczyk (1986).  $\hat{\mathbf{F}} = \frac{1}{\sqrt{T}} \{\psi_1, \psi_2, \dots, \psi_J\}$ , where  $\{\psi_1, \psi_2, \dots, \psi_J\}$  are eigenvectors corresponding to the first  $J$  eigenvalues of  $\frac{1}{n} \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}$ .

Thus,  $\hat{\mathbf{F}} \rightarrow_p \mathbf{F}$  follows. □

**Proof of Theorem 2.6.2:** Given  $\hat{\mathbf{F}}$ , we have:

$$\hat{\mathbf{G}}(\mathbf{X}) = \hat{\mathbf{Y}} \hat{\mathbf{F}} (\hat{\mathbf{F}}^\top \hat{\mathbf{F}})^{-1},$$

as  $\hat{\mathbf{F}}^\top \hat{\mathbf{F}} = \mathbf{I}_J$ , therefore,

$$\hat{\mathbf{G}}(\mathbf{X}) = \tilde{\mathbf{Y}} \hat{\mathbf{F}}.$$

Then we need to show:

$$E((\hat{\mathbf{G}}(\mathbf{X}_i) - \mathbf{G}(\mathbf{X}_i))^2) = 0.$$

Take the sample analogue,

$$\frac{1}{n} (\hat{\mathbf{G}}(\mathbf{X}) - \mathbf{G}(\mathbf{X}))^\top (\hat{\mathbf{G}}(\mathbf{X}) - \mathbf{G}(\mathbf{X})).$$

Given:

$$\mathbf{G}(\mathbf{X}) = \Phi(\mathbf{X}) \mathbf{B} + \mathbf{R}^\theta(\mathbf{X}).$$

$$\hat{\mathbf{G}}(\mathbf{X}) = (\Phi(\mathbf{X}) \mathbf{B} + \mathbf{P} \Lambda + \mathbf{P} \mathbf{R}^\theta(\mathbf{X})) \mathbf{F}^\top \hat{\mathbf{F}} + \mathbf{P} \mathbf{U} \mathbf{D}_T \hat{\mathbf{F}}$$

Furthermore,

$$\mathbf{G}(\mathbf{X}) - \hat{\mathbf{G}}(\mathbf{X}) = (\Phi(\mathbf{X}) \mathbf{B} + \mathbf{P} \Lambda + \mathbf{P} \mathbf{R}^\theta(\mathbf{X})) \mathbf{F}^\top \hat{\mathbf{F}} + \mathbf{P} \mathbf{U} \mathbf{D}_T \hat{\mathbf{F}} - \Phi(\mathbf{X}) \mathbf{B} - \mathbf{R}^\theta(\mathbf{X}) = \mathbf{q}_1 + \mathbf{q}_2 + \mathbf{q}_3 + \mathbf{q}_4.$$

Similar to the Proof of Theorem 2.6.1,

$$\frac{1}{n} (\hat{\mathbf{G}}(\mathbf{X}) - \mathbf{G}(\mathbf{X}))^\top (\hat{\mathbf{G}}(\mathbf{X}) - \mathbf{G}(\mathbf{X})) \rightarrow^P \frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_1 + \frac{1}{n} \mathbf{q}_3^\top \mathbf{q}_3 + \frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_3 + \frac{1}{n} \mathbf{q}_3^\top \mathbf{q}_1.$$

For the first term,

$$\frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_1 = \hat{\mathbf{F}}^\top \mathbf{F} (\Phi(\mathbf{X}) \mathbf{B} + \mathbf{P} \Lambda + \mathbf{P} \mathbf{R}^\theta(\mathbf{X}))^\top (\Phi(\mathbf{X}) \mathbf{B} + \mathbf{P} \Lambda + \mathbf{P} \mathbf{R}^\theta(\mathbf{X})) \mathbf{F}^\top \hat{\mathbf{F}},$$



due to

$$\frac{1}{n} \sum_{j=1}^3 \mathbf{e}_2^\top \mathbf{e}_j + \frac{1}{n} \sum_{j=1}^3 \mathbf{e}_j^\top \mathbf{e}_2 \rightarrow^P \mathbf{0},$$

and

$$\frac{1}{n} \mathbf{e}_1^\top \mathbf{e}_1 \rightarrow^P \mathbf{F} \frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n} \mathbf{F}^\top$$

then,

$$\frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_1 \rightarrow^P \hat{\mathbf{F}}^\top \mathbf{F} \frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n} \mathbf{F}^\top \hat{\mathbf{F}}.$$

Theorem 2.6.1 and Assumption 2.6.1 give  $\hat{\mathbf{F}} \rightarrow \mathbf{F}$  and  $\mathbf{F}^\top \mathbf{F} = \mathbf{I}_J$ , therefore:

$$\frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_1 \rightarrow^P \frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n},$$

Similarly,

$$\begin{aligned} \frac{1}{n} \mathbf{q}_3^\top \mathbf{q}_3 &\rightarrow^P \frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n}, \\ \frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_3 &\rightarrow^P -\frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n}, \\ \frac{1}{n} \mathbf{q}_3^\top \mathbf{q}_1 &\rightarrow^P -\frac{\mathbf{B}^\top \Phi^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{B}}{n}. \end{aligned}$$

Therefore,

$$\frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_1 + \frac{1}{n} \mathbf{q}_3^\top \mathbf{q}_3 + \frac{1}{n} \mathbf{q}_1^\top \mathbf{q}_3 + \frac{1}{n} \mathbf{q}_3^\top \mathbf{q}_1 \rightarrow 0.$$

Then,

$$\frac{1}{n} (\hat{\mathbf{G}}(\mathbf{X}) - \mathbf{G}(\mathbf{X}))^\top (\hat{\mathbf{G}}(\mathbf{X}) - \mathbf{G}(\mathbf{X})) \rightarrow^P 0,$$

thus,

$$\hat{\mathbf{G}}(\mathbf{X}) \rightarrow^P \mathbf{G}(\mathbf{X}).$$

Then Theorem 2.6.2 follows. □

**Proof of Theorem 2.6.3 :** Let  $\dot{\mathbf{Y}} = \frac{1}{T} (\mathbf{Y} - \hat{\mathbf{G}}(\mathbf{X}) \hat{\mathbf{F}}) \mathbf{1}_T$ . By substituting the restriction, we have the Lagrangian equation:

$$\min_A (\dot{\mathbf{Y}} - \Phi(\mathbf{X}) \mathbf{A})^\top (\dot{\mathbf{Y}} - \Phi(\mathbf{X}) \mathbf{A}) + \lambda \hat{\mathbf{G}}^\top(\mathbf{X}) \Phi(\mathbf{X}) \mathbf{A} \quad (\text{A.3})$$

Then we take the first order condition with respect to  $\mathbf{A}$  and  $\boldsymbol{\lambda}$  separately, and we obtain:

$$\begin{pmatrix} 2\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}) & \boldsymbol{\Phi}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}) \\ \hat{\mathbf{G}}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X})^\top & \mathbf{0} \end{pmatrix} \begin{pmatrix} \hat{\mathbf{A}} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} 2\boldsymbol{\Phi}(\mathbf{X})^\top \dot{\mathbf{Y}} \\ \mathbf{0} \end{pmatrix}. \quad (\text{A.4})$$

Under Assumption 2.6.1, the above matrices are invertible, which can be written as:

$$\begin{pmatrix} \hat{\mathbf{A}} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} 2\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}) & \boldsymbol{\Phi}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}) \\ \hat{\mathbf{G}}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X})^\top & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} 2\boldsymbol{\Phi}(\mathbf{X})^\top \dot{\mathbf{Y}} \\ \mathbf{0} \end{pmatrix}. \quad (\text{A.5})$$

Therefore, we obtain:

$$\hat{\mathbf{A}} = \mathbf{Q}\tilde{\mathbf{A}},$$

where

$$\mathbf{Q} = \mathbf{I} - (\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}))^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}) (\hat{\mathbf{G}}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}))^{-1} \hat{\mathbf{G}}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}),$$

$$\tilde{\mathbf{A}} = \frac{1}{T} (\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}))^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top \dot{\mathbf{Y}} \mathbf{1}_T.$$

Furthermore, let  $\boldsymbol{\Xi} = \boldsymbol{\Phi}(\mathbf{X})\hat{\mathbf{A}} - \mathbf{h}(\mathbf{X}) = \boldsymbol{\Phi}(\mathbf{X})\mathbf{Q}\tilde{\mathbf{A}} - \boldsymbol{\Phi}(\mathbf{X})\mathbf{A} - \mathbf{R}^\mu(\mathbf{X})$ .

Under the restriction  $\hat{\mathbf{G}}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X})\mathbf{A} = \mathbf{0}$ , we can obtain:

$$\boldsymbol{\Xi} = \boldsymbol{\Phi}(\mathbf{X})\mathbf{M}(\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}))^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top \frac{1}{T} (\boldsymbol{\Phi}(\mathbf{X})\mathbf{A} + \mathbf{R}^\mu(\mathbf{X}) + \Gamma + (\Lambda + \mathbf{R}^\theta(\mathbf{X}))\mathbf{F}') \mathbf{1}_T - \boldsymbol{\Phi}(\mathbf{X})\mathbf{A} - \mathbf{R}^\mu(\mathbf{X}). \quad (\text{A.6})$$

Furthermore, we have:

$$\boldsymbol{\Phi}(\mathbf{X})\mathbf{M}(\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}))^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top = (\mathbf{I} - \boldsymbol{\Phi}(\mathbf{X})(\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}))^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}) (\hat{\mathbf{G}}(\mathbf{X})^\top \hat{\mathbf{G}}(\mathbf{X}))^{-1} \hat{\mathbf{G}}(\mathbf{X})^\top) \mathbf{P}. \quad (\text{A.7})$$

And then, substitute Equation A.7 into Equation A.6 and under Assumption 2.2.1 and Theorem 2.6.2:

$$\boldsymbol{\Xi} = \boldsymbol{\Phi}(\mathbf{X})\mathbf{A} - \boldsymbol{\Phi}(\mathbf{X})\mathbf{A} - \mathbf{R}^\mu(\mathbf{X}).$$

$$\mathbf{R}^\mu(\mathbf{X}) \rightarrow \mathbf{0} \text{ as } n \rightarrow \infty,$$

therefore,

$$\frac{1}{n} \boldsymbol{\Xi}^\top \boldsymbol{\Xi} \rightarrow \mathbf{0}.$$

And the Theorem 2.6.3 follows.  $\square$

**Proof of Theorem 2.6.4 :** Define  $Z = \max_{\{1 \leq p \leq P, 1 \leq h \leq H_n\}} \{|\hat{\alpha}_{ph}|/\hat{\sigma}_{ph}\}$ . Under Assumption 3, we have

$$\hat{\alpha}_{ph}/\hat{\sigma}_{ph} | \mathbf{H}_0 \rightarrow^d N(0, 1).$$

Therefore, under the  $\mathbf{H}_0$ , we have:

$$\begin{aligned} e^{tE(Z)} &\leq E[e^{tZ}] \\ &= E[\max\{t|\hat{\alpha}_{ph}|/\hat{\sigma}_{ph}\}] \\ &\leq \sum_{p=1, h=1}^{p=P, h=H_n} E[e^{t|\hat{\alpha}_{ph}|/\hat{\sigma}_{ph}}] \\ &= PH_n e^{t^2/2}. \end{aligned}$$

Then take the logarithm of both sides we can obtain:

$$E[Z] \leq \frac{\log PH_n}{t} + \frac{t}{2}.$$

If we set  $t = \sqrt{2 \log PH_n}$  to minimise  $\frac{\log PH_n}{t} + \frac{t}{2}$ , then we have:

$$E[Z] \leq \sqrt{2 \log PH_n}.$$

Therefore, we can bound the  $|\hat{\alpha}_{ph}|/\hat{\sigma}_{ph}$  by  $\sqrt{2 \log PH_n}$ . □

**Proof of Theorem 2.6.5 :** To proof

$$\Pr(\text{reject } H_0 | \hat{\mathcal{M}} \neq \emptyset) \rightarrow 1,$$

equivalently, we need to prove

$$\Pr(S_0 + S_1 | \hat{\mathcal{M}} \neq \emptyset) \rightarrow 1$$

$S_0 = H_n \sum_{p=1}^P \mathbf{I}(\sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n)$ , as  $H_n = n^v \rightarrow \infty$  when  $n \rightarrow \infty$ .

Once  $\hat{\mathcal{M}} \neq \emptyset$ , then  $\sum_{p=1}^P \mathbf{I}(\sum_{h=1}^{H_n} |\hat{\alpha}_{ph}|/\hat{\sigma}_{ph} \geq \eta_n) \geq 1$ , therefore,  $S_0 \rightarrow \infty$  as  $n \rightarrow \infty$ . Meanwhile  $F_q = O(1)$ , we can show that:

$$\Pr(S_0 + S_1 > F_q | \hat{\mathcal{M}} \neq \emptyset) \rightarrow 1.$$

Then the Theorem 2.6.5 follows. □

### A.1.3 Proofs of Chapter 3

**Proof of Theorem 3.4.1 :** Write the subportfolio vector  $\mathbf{Q}_t(\mathbf{X})$  as:

$$\mathbf{Q}_t(\mathbf{X}) = \sum_{i=1}^n \mathbf{B}(\mathbf{X}_i) y_{it}.$$

Because:

$$y_{it} = G(\mathbf{X}_i) \mathbf{F}_t + \varepsilon_{it}.$$

Then, substitute  $y_{it}$  into  $\mathbf{Q}_t(\mathbf{X})$ :

$$\begin{aligned} \mathbf{Q}_t(\mathbf{X}) &= \sum_{i=1}^n \mathbf{B}(\mathbf{X}_i) (G(\mathbf{X}_i) \mathbf{F}_t + \varepsilon_{it}) \\ &= \frac{1}{n} \sum_{i=1}^n \Gamma G(\mathbf{X}_i)^\top (G(\mathbf{X}_i) \mathbf{F}_t + \varepsilon_{it}) \\ &= \Gamma \left( \frac{1}{n} \sum_{i=1}^n (G(\mathbf{X}_i)^\top G(\mathbf{X}_i)) \right) \mathbf{F}_t + \frac{1}{n} \sum_{i=1}^n (G(\mathbf{X}_i)^\top \varepsilon_{it}) \end{aligned}$$

Given:

$$p \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n G(\mathbf{X}_i)^\top G(\mathbf{X}_i) = \mathbf{M}^G,$$

where  $\mathbf{M}^G$  is an identity matrix, and

$$E(\varepsilon_{it} | \mathbf{X}_i, \mathbf{F}_t) = 0.$$

Therefore, we have:

$$p \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n G(\mathbf{X}_i)^\top \varepsilon_{it} = 0.$$

Thus,

$$p \lim_{n \rightarrow \infty} \mathbf{Q}_t(\mathbf{X}) = \Gamma \mathbf{M}^G \mathbf{F}_t = \Gamma \mathbf{F}_t.$$

This shows that the factor-mimicking portfolio is a linear combination of risk factors given  $\Gamma$  is a non-zero matrix.  $\square$

**Proof of Theorem 3.4.2 :** Let  $\tilde{\mathbf{F}}$  represent the demeaned risk factor matrix while

$$\tilde{\mathbf{y}}_t = G(\mathbf{X}) \tilde{\mathbf{F}}_t + \varepsilon_t.$$

Correspondingly, we have:

$$\begin{aligned}\tilde{\mathbf{Q}}_t(\mathbf{X}) &= \mathbf{B}(\mathbf{X}_i)\tilde{\mathbf{y}}_t \\ &= \frac{1}{n}\Gamma\mathbf{G}(\mathbf{X})^\top\tilde{\mathbf{y}}_t.\end{aligned}$$

And then,

$$\begin{aligned}E(\tilde{\mathbf{Q}}_t(\mathbf{X})\tilde{\mathbf{Q}}_t(\mathbf{X})^\top|\mathbf{X}) &= \Gamma\left(\frac{1}{n}\mathbf{G}(\mathbf{X})^\top\mathbf{G}(\mathbf{X})\right)E(\tilde{\mathbf{F}}\tilde{\mathbf{F}}^\top)\left(\frac{1}{n}\mathbf{G}(\mathbf{X})^\top\mathbf{G}(\mathbf{X})\right)\Gamma^\top + \\ &\quad \Gamma\left(\frac{1}{n}\mathbf{G}(\mathbf{X})^\top\frac{1}{n}E(\boldsymbol{\varepsilon}_t\boldsymbol{\varepsilon}_t^\top)\right)\mathbf{G}(\mathbf{X})\Gamma^\top.\end{aligned}$$

Given  $E(\boldsymbol{\varepsilon}_{it}|\mathbf{X}_i, \mathbf{F}_t) = 0$ , the cross terms are  $E(\tilde{\mathbf{F}}\tilde{\boldsymbol{\varepsilon}}_t) = 0$ .

Taking the second term and using the Euclidian matrix norm:

$$\begin{aligned}\|\Gamma\left(\frac{1}{n}\mathbf{G}(\mathbf{X})^\top\frac{1}{n}E(\tilde{\boldsymbol{\varepsilon}}_t\tilde{\boldsymbol{\varepsilon}}_t^\top)\mathbf{G}(\mathbf{X})\right)\Gamma^\top\| &\leq \\ \frac{1}{n}\|\Gamma\frac{1}{n}(\mathbf{G}(\mathbf{X})^\top\mathbf{G}(\mathbf{X}))\Gamma^\top\| \times \|E(\boldsymbol{\varepsilon}_t\boldsymbol{\varepsilon}_t^\top)\| &\rightarrow_{n\rightarrow\infty} \\ \frac{1}{n}\|\Gamma\Gamma^\top\| \times \|E(\boldsymbol{\varepsilon}_t\boldsymbol{\varepsilon}_t^\top)\| &\rightarrow_{n\rightarrow\infty} 0\end{aligned}$$

The conclusion of the above formula is due to

$$p \lim_{n\rightarrow\infty} \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\mathbf{X}_i)^\top \mathbf{G}(\mathbf{X}_i) = \mathbf{M}^G,$$

and  $\|E(\boldsymbol{\varepsilon}_t\boldsymbol{\varepsilon}_t^\top)\|$  has bounded eigenvalues for all  $n$ .

Furthermore, the well-chosen coefficient matrix  $\Gamma$  can give:

$$E(\tilde{\mathbf{Q}}_t(\mathbf{X})\tilde{\mathbf{Q}}_t(\mathbf{X})^\top|\mathbf{X}) = \mathbf{I}_{JJ}$$

□

**Proof of Theorem 3.4.3 :** We decompose the investment returns of optimal asset-by-asset portfolio and risk-free rate as:

$$r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^*,$$

where the  $\mathbf{R}_F$  is the optimal factor returns since the return generation function states the risk premiums come from risk factors. The  $\boldsymbol{\varepsilon}_t^*$  is the zero mean idiosyncratic returns.

Since

$$E(\boldsymbol{\varepsilon}_{it}|\mathbf{X}_i, \mathbf{F}_t) = 0,$$

it follows from the second-order stochastic dominance that the expected utility has the following relationship:

$$E(u(r_{ft} + \mathbf{R}_F)) > E(u(r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^*)).$$

Because zero mean  $\boldsymbol{\varepsilon}_t^*$  only contributes variance rather than returns.

According to Theorem 3.4.1 and Equation 3.18, the restricted portfolio optimally combines the factors' returns. Therefore, our two-stage portfolio's return can be written as:

$$r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^{**},$$

where the only difference is the idiosyncratic returns. Our goal now is to show that:

$$E(u(r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^{**})|\mathbf{X}, \mathbf{z}_t) \rightarrow^{n \rightarrow \infty} E(u(r_{ft} + \mathbf{R}_F)|\mathbf{X}, \mathbf{z}_t).$$

Next, we take the Taylor expansion of  $u(r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^{**})$  around  $r_{ft} + \mathbf{R}_F$ :

$$u(r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^{**}) = u(r_{ft} + \mathbf{R}_F) + \frac{d}{d(r_{ft} + \mathbf{R}_F)} u(r_{ft} + \mathbf{R}_F) \boldsymbol{\varepsilon}_t^{**} + \frac{d^2}{d(r_{ft} + \mathbf{R}_F)^2} u(r_{ft} + \mathbf{R}_F)^2 \boldsymbol{\varepsilon}_t^{**2}.$$

We take the expectation on both sides, given  $E(\boldsymbol{\varepsilon}_t^{**}) = 0$  and  $\frac{d^2 u(\cdot)}{dW^2} \geq -c$ . Therefore, we have:

$$E(u(r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^{**})) \geq E(u(r_{ft} + \mathbf{R}_F)) - cE(\boldsymbol{\varepsilon}_t^{**2}),$$

where  $p \lim_{n \rightarrow \infty} E(\boldsymbol{\varepsilon}_t^{**2}) = 0$ , according to Theorem 3.4.2. Therefore, we have :

$$p \lim_{n \rightarrow \infty} E(u(r_{ft} + \mathbf{R}_F + \boldsymbol{\varepsilon}_t^{**})|\mathbf{X}, \mathbf{z}_t) - E(u(r_{ft} + \mathbf{R}_F)|\mathbf{X}, \mathbf{z}_t) = 0,$$

which completes the proof. □

# Appendix B

## Tables and Figures

### B.1 Characteristic Description

Table B.1 Characteristic Details

<b>Name</b>	<b>Description</b>	<b>Reference</b>
A2ME	We define assets-market cap as total assets (AT) over market capitalization as of December t-1. Market capitalization is the product of shares outstanding (SHROUT) and price(PRC).	Bhandari (1988)
AT	Total assets (AT)	Gandhi and Lusting (2015)
ATO	Net sales over lagged net operating assets. Net operating assets are the difference between operating assets and operating liabilities. Operating assets are total assets (AT) minus cash and short-term investments (CHE), minus investment and other advances (IVAO). Operating liabilities are total assets (AT), minus debt in current liabilities(DLC),minus long-term debt (DLTT),minus minority interest (MIB), minus preferred stock (PSTK), minus common equity (CEQ).	Soliman(2008)

BEME	Ratio of book value of equity to market value of equity. Book equity is shareholder equity (SH) plus deferred taxes and investment tax credit (TXDITC), minus preferred stock (PS). SH is shareholder's equity (SEQ). If missing, SH is the sum of common equity (CEQ) and preferred stock (PS). If missing, SH is the difference between total assets (AT) and total liabilities (LT). Depending on availability, we use the redemption (item PSTKRV), liquidating (item PSTKL), or par value (item PSTK) for PS. The market value of equity is as of December t-1. The market value of equity is the product of shares outstanding (SHROUT) and price (PRC).	Rosenberg, Reid and Lanstein (1985) Davis, Fama, and French (2000)
C	Ration of cash and short-term investments (CHE) to total assets (AT)	Palazzo
C2D	Cash flow to price is the ratio of income and extraordinary items (IB) and depreciation and amortization (dp) to total liabilities (LT).	
CTO	We define caoital turnover as ratio of net sales (SALE) to lagged total assets (AT).	Haugen and Baker (1996)
Debt2P	Debt to price is the radio of long-term debt (DLTT) and debt in current liabilities (DLC) to the market capitalization as of December t-1 . Market capitalization is the product of shares outstanding (SHROUT) and price (PRC).	Litzenberger and Ramaswamy (1979)
$\Delta ceq$	The percentage change in the book value of equity (CEQ).	Richardson et al. (2005)
$\Delta(\Delta Gm - Sales)$	The difference in the percentage change in gross margin and the percentage change in sales (SALE). We define gross margin as the difference in sales (SALE) and costs of goods sold (COGS).	Abarbanell and Bushee (1997)
$\Delta ShROUT$	The definition of the percentage change in shares outstanding (SHROUT).	Pontiff and Woodgate (2008)



$\Delta PI2A$	We define the change in property, plants ,and equip- ment as changes in property,plants,and equipment (PPEGT) and inventory (INVT) over lagged total assets (TA).	Lyandres , Sun, and Zhang (2008)
DTO	We define turnover as ratio of daily volume (VOL) to shares outstanding (SHROUT) minus the daily market turnover and de-trend it by its 180 trading day median. We scale down the volume of NAS- DAQ securities by 38% after 1997 and by 50% before that to address the issue of double-counting of volume for NASDAQ securities.	Garfinkel (2009); Anderson and Dyl (2005)
E2P	We define earnings to price as the ratio of income before extraordinary items (IB) to the market cap- italization as December t-1 Market capitalization is the product of share outstanding (SHROUT) and price (PRC).	Basu (1983)
EPS	We define earnings per share as the ratio of income before extraordinary items (IB) to share outstanding (SHROUT) as of December t-1	Basu (1997)
Investment	We define investment as the percentage year-on- year growth rate in total assets (AT).	Cooper, Gulen and Schill(2008)
IPM	We define pre-tax profit margin as ratio of pre-tax income (PI) to sales (SALE).	
Lev	leverage is the ratio of long-term debt (DLTT) and debt in the current liabilities (DLC) to the sum of long-term debt, debt in current liabilities, and stock- holders' equity (SEQ)	Lewenllen (2015)
LME	Size is the total market capitalization of the pre- vious month defined as price (PRC) times shares outstanding (SHROUT)	Fama and French (1992)
Turnover	Turnover is last month's volume (VOL) over shares outstanding (SHROUT).	Datar, Naik and Radcliffe (1998)
PCM	The price-to-cost margin is the difference between net sales (SALE) and costs of goods sold (COGS) divided by net sales (SALE).	Gorodnichenko and Weber (2016) and D'Acunto, Liu, Pflueger and Wcber (2017)

PM	The profit margin is operating income after depreciation (OIADP) over sales (SALE)	Soliman (2008)
Q	Tobin's Q is total assets (AT), the market value of equity (SHROUT times PRC) minus cash and short-term investments (CEQ) minus deferred taxes (TXDB) scaled by total assets (AT).	
ROA	Return-on-assets is income before extraordinary items (IB) to lagged total assets (AT).	Balakrishnan, Bartov and Faurel (2010)
ROC	ROC is the ratio of market value of equity (ME) plus long-term debt (DLTT) minus total assets to Cash and Short-Term Investments (CHE).	Chandrashekar and Rao (2009)
ROE	Return-on-equity is income before extraordinary items (IB) to lagged book-value of equity.	in Haugen and Baker (1996)
$r_{12-2}$	We define momentum as cumulative return from 12 months before the return prediction to two months before.	Fama and French (1996)
$r_{12-7}$	We define intermediate momentum as cumulative return from 12 months before the return prediction to seven months before.	Novy-Marx (2012)
$r_{6-2}$	We define $r_{6-2}$ as cumulative return from 6 months before the return prediction to two months before.	Jegadeesh and Titman (1993)
$r_{2-1}$	We define short-term reversal as lagged one-month return.	Jegadeesh(1990)
S2C	Sales-to-cash is the ratio of net sales (SALE) to Cash and Short-Term Investments (CHE).	following Ou and Penman (1989)
Sales-G	Sales growth is the percentage growth rate in annual sales (SALE).	Lakonishok, Shleifer, and Vishmy (1994)
SGA2S	SGA to sales is the ratio of selling, general and administrative expenses (XSGA) to net sales (SALE).	

## B.2 Figures and Tables

### B.2.1 Figures of Chapter 2

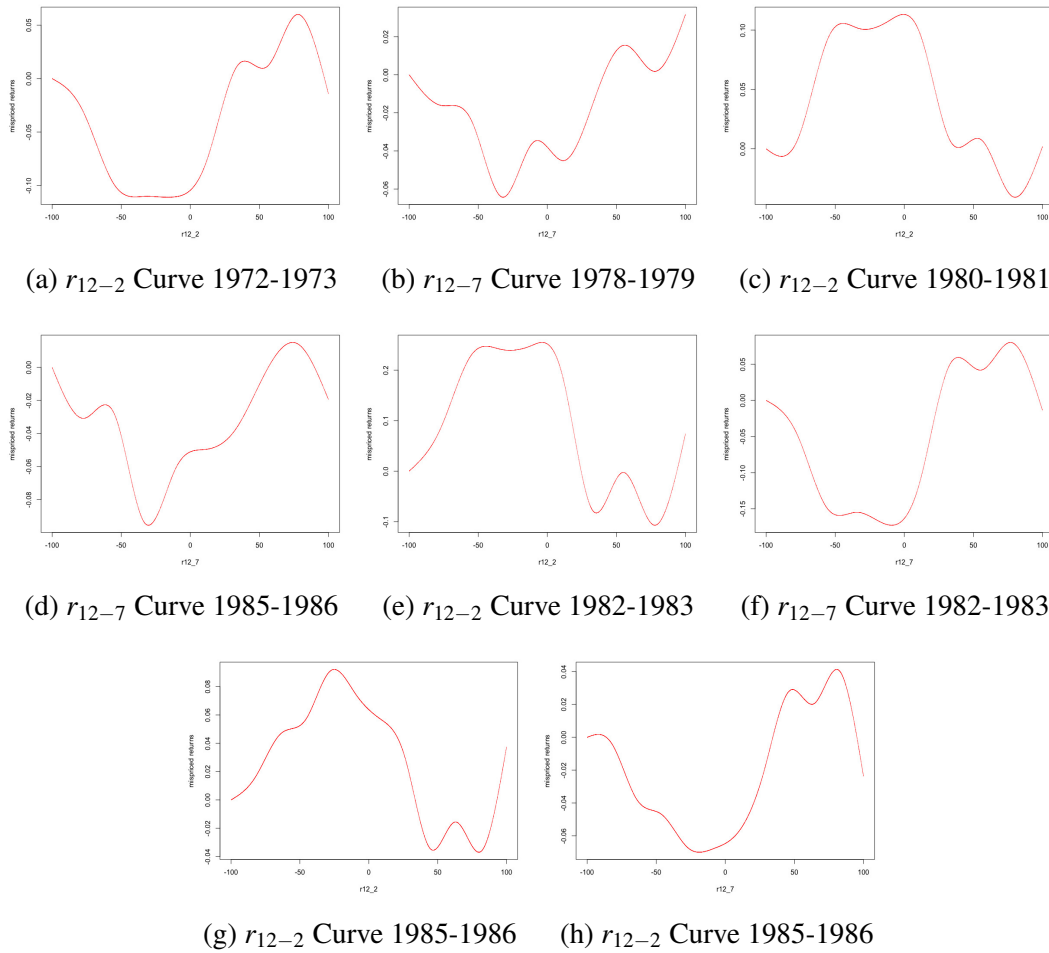


Figure B.1 Mispricing Characteristic Curve of standardized  $r_{12-2}$  and  $r_{12-7}$

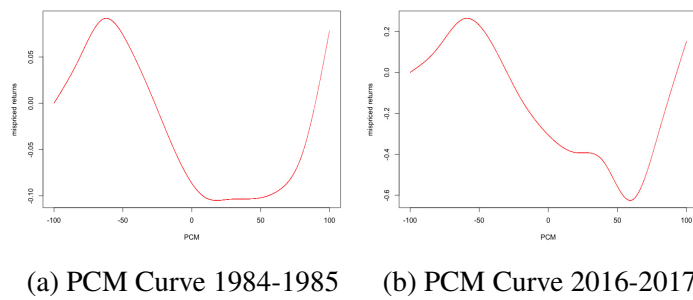


Figure B.2 Mispricing Characteristic Curve of standardized PCM

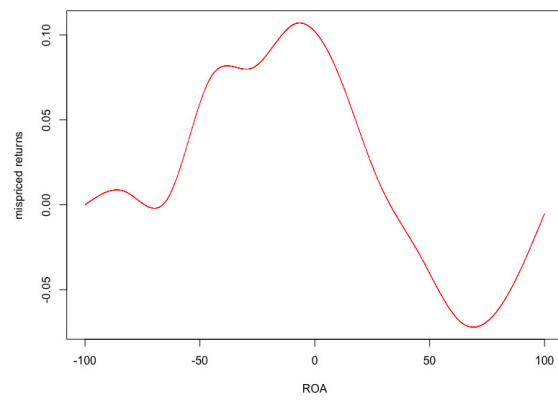
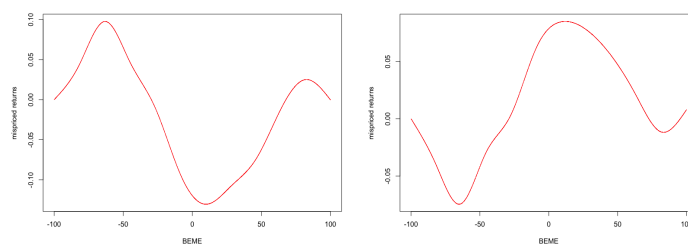
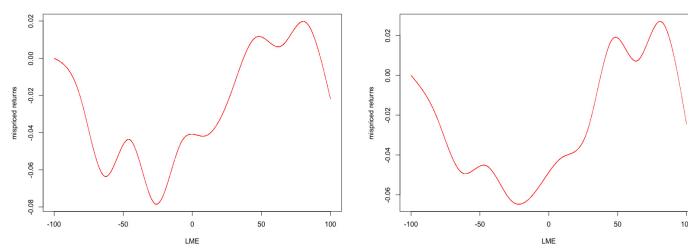


Figure B.3 Mispricing Characteristic Curve of standardized ROA in 1988-1989



(a) BEME Curve 1995-1996 (b) BEME Curve 1996-1997

Figure B.4 Mispricing Characteristic Curve of standardized BEME



(a) LME Curve 1998-1999 (b) LME Curve 2000-2001

Figure B.5 Mispricing Characteristic Curve of standardized LME

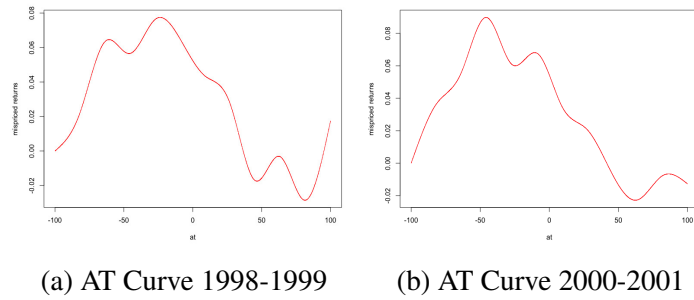


Figure B.6 Mispricing Characteristic Curve of standardized AT

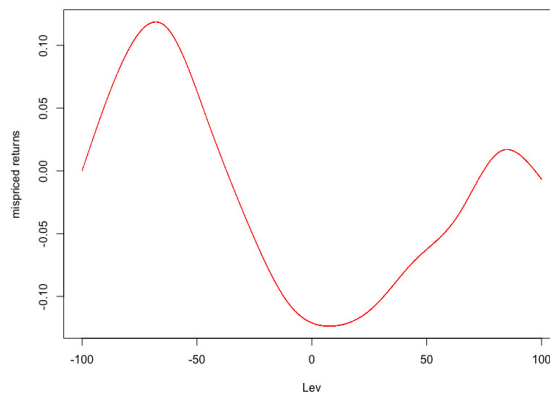


Figure B.7 Mispricing Characteristic Curve of standardized LEV in 2002-2003

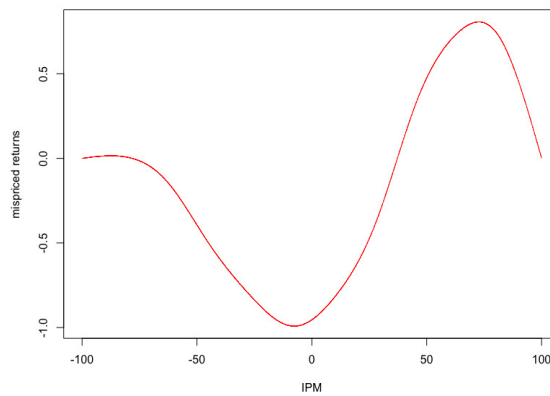


Figure B.8 Mispricing Characteristic Curve of standardized IPM in 2004-2005

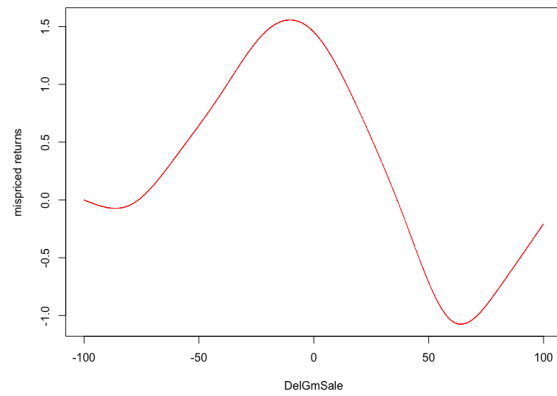


Figure B.9 Mispricing Characteristic Curve of standardized DelGmSale in 2015-2016

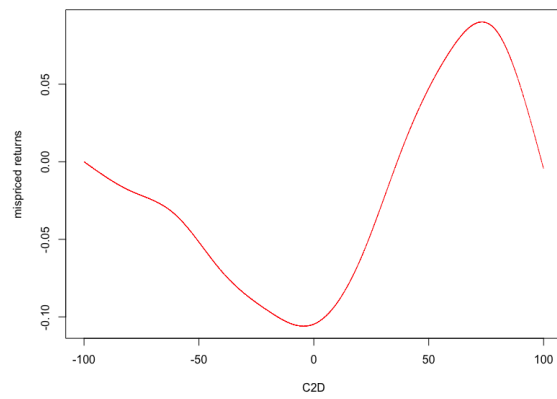


Figure B.10 Mispricing Characteristic Curve of standardized C2D in 2016-2017

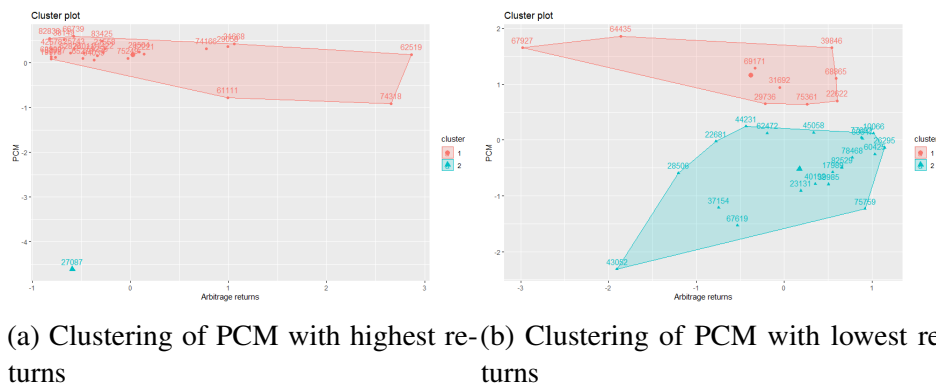
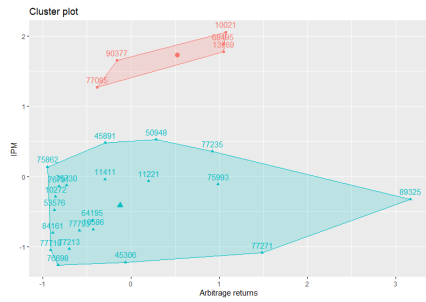
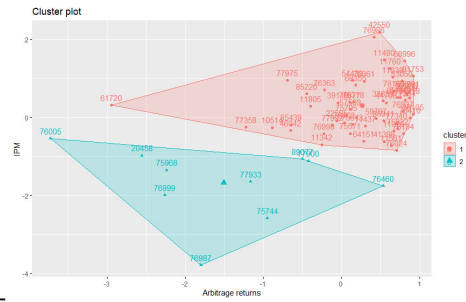


Figure B.11 Clustering of PCM 1986-1987



(a) Clustering of IPM with highest returns



(b) Clustering of IPM with lowest returns

Figure B.12 Clustering of IPM 2004-2005

### B.2.2 Tables of Chapter 3

Table B.2 Annual Correlation Between Subportfolios and Risk Factors 1-20

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
$f_1$	0.61	-0.49	0.68	0.69	0.66	0.73	-0.10	-0.16	0.72	0.79	0.73	0.69	0.77	-0.17	0.64	0.87	0.65	0.71	0.51	0.75
$f_2$	0.89	0.64	0.70	0.75	0.65	0.79	0.60	0.91	0.71	0.79	0.90	0.92	0.97	-0.01	0.79	0.74	0.79	0.86	0.88	0.41
$f_3$	0.80	0.81	0.63	0.40	0.66	0.39	0.75	0.64	0.40	0.72	0.86	0.74	0.96	0.56	0.72	0.86	0.84	0.66	0.78	0.55

This table shows the annual correlation between factor-mimicking subportfolios and corresponding risk factors from Jul. 1967- Jun.1987.

Table B.3 Annual Correlation Between Subportfolios and Risk Factors 21-40

	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40
$f_1$	0.78	0.86	0.54	0.58	0.22	0.76	0.68	0.16	0.60	0.74	0.75	0.67	0.72	0.73	0.69	0.67	0.96	0.66	0.73	0.76
$f_2$	0.77	0.79	0.73	-0.24	0.93	0.66	0.86	-0.06	0.87	0.81	0.94	0.76	0.22	-0.06	0.97	0.94	0.92	0.36	0.93	0.93
$f_3$	-0.40	0.72	0.76	0.64	0.71	0.83	0.70	0.57	0.73	0.67	0.72	0.58	0.66	0.60	0.79	0.77	0.92	0.68	0.79	0.72

This table shows the annual correlation between factor-mimicking subportfolios and corresponding risk factors from Jul. 1987- Jun.2007.

Table B.4 Annual Correlation Between Subportfolios and Risk Factors 41-50

	41	42	43	44	45	46	47	48	49	50
$f_1$	-0.22	0.73	-0.15	0.72	-0.55	-0.08	0.33	0.74	0.63	0.77
$f_2$	0.63	0.76	0.29	0.73	0.65	0.34	-0.23	0.81	0.83	0.81
$f_3$	0.64	0.58	0.65	0.61	0.82	0.51	0.87	0.86	0.82	0.53

This table shows the annual correlation between factor-mimicking subportfolios and corresponding risk factors from Jul. 2007- Jun.2017.