

UNIVERSITÀ DEGLI STUDI DI PADOVA  
DIPARTIMENTO DI SCIENZE STATISTICHE  
CORSO DI LAUREA MAGISTRALE IN  
SCIENZE STATISTICHE



**Revisiting the Small Cap Effect:  
an empirical analysis using cross-sectional factors**

Relatore Prof. Massimiliano Caporin  
Dipartimento di Scienze Statistiche

Laureando Giorgio Albanese  
Matricola 2082884

Anno Accademico 2023/2024



# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Literature Review (A)</b>	<b>5</b>
1.1 Asset Pricing Models . . . . .	5
1.1.1 The Capital Asset Pricing Model (CAPM) . . . . .	5
1.1.2 Arbitrage Pricing Theory (APT) . . . . .	6
1.2 Time-Series Factors . . . . .	7
1.2.1 Fama-French Three and Five Factor Models . . . . .	7
1.2.2 Production of Time Series Factors (Fama and French, 2023) . . . . .	8
1.3 Cross-Sectional Factors . . . . .	8
1.4 Two Pass Procedure . . . . .	10
1.5 Factor Strength . . . . .	11
<b>2 Literature Review (B)</b>	<b>17</b>
2.1 Empirical Results of the Fama-French Five-Factor Model . . . . .	17
2.1.1 Focus on the SMB Factor . . . . .	17
2.1.2 Summary . . . . .	18
2.2 Empirical evidence for Time-Series Models using Cross-Sectional Factors (Fama and French, 2020) . . . . .	19
2.2.1 The framework . . . . .	19
2.2.2 Comparisons and summary statistics . . . . .	19
2.3 Factor Strength (Pesaran and Smith, 2021) . . . . .	22
2.3.1 Data and Methodology . . . . .	22
2.3.2 Results: Factor Strength Estimation . . . . .	22
<b>3 Data</b>	<b>27</b>
3.1 Data collection . . . . .	27
3.2 Pre-Processing . . . . .	29
3.3 Data coverage . . . . .	31
<b>4 Cross-Sectional factors</b>	<b>41</b>
4.1 The framework . . . . .	41
4.2 The baseline estimate . . . . .	42
4.3 Challenging the Small Cap Effect . . . . .	52
4.3.1 Including other variables . . . . .	52

---

4.3.2	The role of Megacaps . . . . .	65
4.3.3	Other drivers of the size factor . . . . .	73
<b>5</b>	<b>Time-Series analysis</b>	<b>75</b>
5.1	Time-Series results . . . . .	75
5.2	Evaluating Factor Strength . . . . .	82
<b>6</b>	<b>Assessing Methodological Weaknesses and Robustness</b>	<b>89</b>
6.1	On the robustness of the treatment of outliers . . . . .	89
6.2	On the omitted variable bias . . . . .	92
	<b>Conclusion</b>	<b>93</b>
	<b>Appendix</b>	<b>97</b>
	<b>Bibliography</b>	<b>107</b>





# Introduction

The focus of this thesis is the empirical analysis of the *small cap effect*, a widely observed phenomenon in financial markets where small-cap stocks often exhibit higher average returns compared to their larger counterparts. This size premium, or small-cap bias, has traditionally been seen as value-additive in equity markets, contributing positively to diversified portfolios. However, recent years have witnessed a reversal of this trend, with large-cap stocks frequently outperforming smaller firms, raising questions about the persistence and underlying drivers of the small-cap effect.

A primary objective of this research is to investigate whether the structural factors that originally led to the small-cap premium are still active in today's markets, or whether this effect has diminished in recent years. Additionally, this thesis seeks to understand the role of fundamental drivers in shaping the dynamics of the size factor, analyzing the extent to which firm characteristics and market conditions contribute to its behavior over time.

Historically, numerous studies have explored the reasons behind the superior returns of small-cap stocks, attributing this phenomenon to factors such as higher perceived risk, lower liquidity, or unique firm characteristics. Pioneering research by Fama and French (1993) integrated size as a core component of their three- and five-factor models, constructing the size factor by measuring the return differentials between portfolios created ex ante through iterative sorts across relevant characteristics. However, while this methodology has become foundational, it also has notable limitations, particularly in terms of controlling for interactions between multiple characteristics. When characteristics such as profitability or value are interdependent with size, simple portfolio sorting may be insufficient to isolate the pure effect of size, as multiple sorting intersections quickly become infeasible and may introduce biases from omitted variables.

In response, this thesis proposes an alternative approach to factor construction, diverging from the traditional Fama-French methodology. Specifically, we implement a

cross-sectional regression-based framework, where factors are constructed by conducting multiple cross-sectional regressions at each point in time. The resulting factors are then stacked into time series, generating an alternative set of factors that are optimized based on the observed returns and firm characteristics in each cross section. This approach provides two primary advantages: first, it enables the retrieval of “purer” factors by allowing for the inclusion of control variables that mitigate the influence of unrelated characteristics; and second, it enhances the explanatory power of the factors giving that they are retrieved through an optimization procedure.

Furthermore, this study addresses some of the methodological critiques commonly associated with cross-sectional factor construction, acknowledging the potential for subjectivity and biases in factor design decisions. Following established protocols from the literature, we critically evaluate the robustness of our framework, with an emphasis on understanding how methodological choices impact the final estimates of factor strength and pervasiveness.

Our dataset includes individual stock data from a range of countries and markets, facilitating a comparative analysis across both geographical and sectoral dimensions. This comprehensive dataset allows us to assess the extent to which the small-cap effect varies across different economic environments, potentially uncovering regional or sectoral differences in the size premium.

To support the validity of our findings, we also examine the strength of the factors retrieved, measuring their influence and persistence over time. This analysis provides insights into the stability of the small-cap effect and offers a foundation for assessing its relevance in contemporary markets.

The structure of this thesis is as follows: Chapters 1 and 2 provide a review of the relevant literature. Chapter 1 focuses on the methodological tools applied in this analysis, detailing various approaches to factor construction and the theoretical frameworks underlying the small-cap effect. Chapter 2 presents an overview of the current understanding of size effects in the literature, reporting some empirical findings. Chapter 3 describes the dataset used, including details on data acquisition, the properties of the data, and the pre-processing methods employed. Chapter 4 presents the cross-sectional factor retrieval process, with an in-depth examination of the impact of other variables on the size factor. Chapter 5 transitions to a time-series analysis, evaluating the time-varying behavior and strength of the factors derived.

In sum, this thesis aims to offer a refined understanding of the small-cap effect by employing a novel approach to factor construction and by situating the analysis within a modern empirical framework. By doing so, we seek to contribute to the ongoing debate on the persistence and drivers of the size premium.



# Chapter 1

## Literature Review (A)

### 1.1 Asset Pricing Models

#### 1.1.1 The Capital Asset Pricing Model (CAPM)

The *Capital Asset Pricing Model (CAPM)*, developed by Sharpe (1964) and Lintner (1965), is considered the precursor to many of the modern asset pricing models. It relates an asset's expected return to its systematic risk, or *beta* ( $\beta$ ), which captures the asset's sensitivity to market movements.

The equation governing the CAPM is given by:

$$R_{it} - R_{ft} = \alpha_i + \beta_i(R_{mt} - R_{ft}) + \epsilon_{it} \quad (1.1)$$

Where:

- $R_{it}$  is the return on asset  $i$  at time  $t$ ,
- $R_{ft}$  is the risk-free rate,
- $R_{mt}$  is the market return,
- $\alpha_i$  represents the asset's return that cannot be attributed to market movements,
- $\beta_i$  represents the systematic risk of the asset,
- $\epsilon_{it}$  is the idiosyncratic component.

This model, while simple and intuitive, does not capture all the complexities of the real world. In particular, empirical research has demonstrated that the CAPM's ability

to explain the cross-section of returns is limited, especially for portfolios of small-cap or high book-to-market ratio stocks.

### 1.1.2 Arbitrage Pricing Theory (APT)

The Arbitrage Pricing Theory (APT) is an alternative to the Capital Asset Pricing Model (CAPM), proposed by Stephen Ross in 1976, to explain the returns of financial assets. APT is based on the idea that the expected return of an asset can be represented as a linear combination of various macroeconomic factors that influence market risk. Unlike CAPM, which considers only a single risk factor (market risk), APT is a multifactor model that includes multiple sources of risk.

The APT model can be formally written as:

$$R_{it} - R_{ft} = \alpha_i + \sum_{k=1}^K \beta_{ki} F_{kt} + \epsilon_{it} \quad (1.2)$$

Where  $F_k$  is the  $k$ th observed factor.

APT does not specify which risk factors should be considered, leaving flexibility for the researcher to identify the most relevant ones. The chosen factors may vary, commonly used factors include: interest rates, inflation and credit spread among others.

One of the main differences between APT and CAPM lies in the flexibility of APT to consider multiple risk factors, while CAPM considers only market risk as the driver of returns. This difference makes APT more adaptable to complex contexts, where multiple macroeconomic factors impact asset returns.

An important aspect to consider when using the APT is the choice of factors. In many empirical applications, the factors used in the APT model are not directly observable macroeconomic variables but rather returns on well-diversified portfolios. This approach has significant implications for both the interpretation and the estimation of the model.

As explained in the book *Rischio e Rendimento* (Pastorello, 2001), using portfolio returns as factors offers practical advantages. First, portfolios can be constructed to proxy for different risk sources, such as size, value, or momentum. By doing so, the factors are directly related to observable investment strategies, making it easier to link the model to real-world applications.

The use of portfolio returns implies that the estimated factor loadings ( $\beta$  coefficients) represent the sensitivity of the asset to these specific portfolios. Moreover, the assumption that these portfolios capture systematic risk effectively is crucial for the validity of the APT model.

## 1.2 Time-Series Factors

### 1.2.1 Fama-French Three and Five Factor Models

To address the shortcomings of the CAPM, the *Three-Factor Model* was introduced (Fama and French, 1992), which supported by empirical evidence incorporates two additional factors beyond the market risk premium:

1. **Size (SMB - Small Minus Big)**: The excess return of small-cap stocks over large-cap stocks.
2. **Value (HML - High Minus Low)**: The excess return of high book-to-market stocks over low book-to-market stocks.

The model becomes:

$$R_{it} - R_{ft} = \alpha_i + b_i(R_{mt} - R_{ft}) + s_iSMB_t + h_iHML_t + \epsilon_{it} \quad (1.3)$$

The Three-Factor Model significantly improved the explanatory power of asset returns, especially in accounting for anomalies observed in small-cap and value stocks. Building on this, later on the *Five-Factor Model* (Fama and French, 2015) has been proposed, adding the following factors:

1. **Profitability (RMW - Robust Minus Weak)**: The excess return of firms with high profitability over those with low profitability.
2. **Investment (CMA - Conservative Minus Aggressive)**: The excess return of firms with conservative investment policies over those with aggressive investment strategies.

The final model can be written as:

$$R_{it} - R_{ft} = \alpha_i + b_i(R_{mt} - R_{ft}) + s_iSMB_t + h_iHML_t + r_iRMW_t + c_iCMA_t + \epsilon_{it} \quad (1.4)$$

The Three-Factor and Five-Factor Models are estimated using *time-series (TS) regressions*, asset returns are regressed on a set of predefined risk factors. A limitation of this approach is that it assumes constant factor loadings over time, which may not always reflect the dynamic nature of asset returns.

### 1.2.2 Production of Time Series Factors (Fama and French, 2023)

The Fama-French factors discussed above are constructed using data from various sources. The primary sources of stock data are the *Center for Research in Security Prices (CRSP)* database and *Compustat*.

The construction of the Fama-French factors begins by sorting stocks from the U.S. market into portfolios based on their characteristics. For instance, the *Size factor* (SMB) is the return difference between small and large capitalization firms. To construct SMB, stocks are split into two groups (small and big) using the median NYSE market capitalization at the end of June each year. This sorting is done annually, and portfolios are reformed at the end of each June, with returns computed starting from July.

Similarly, the *Value factor* (HML) is the return difference between high book-to-market (value) and low book-to-market (growth) stocks. To construct this factor, stocks are first sorted into three groups based on their book-to-market ratio, with the 30th and 70th percentiles of the NYSE market used as thresholds. Each stock is then classified independently from the first sort as small or big based on its size. The intersection of these groups results in six portfolios, and the HML factor is calculated as the average return of the two high book-to-market portfolios minus the two low book-to-market portfolios. The same procedure based on different characteristics is applied to obtain the other factors.

These factors are ubiquitously used in academic research and practical portfolio management. It is important to note that the construction methodology for these factors could change over time, and updates to data (e.g. error corrections) and methodologies could lead to changes in historical factor returns. As a result, the factor returns are updated regularly, reflecting the best available data at each point in time.

## 1.3 Cross-Sectional Factors

As an alternative to the time-series factors discussed above, Fama and French (2020) propose a new set of Factors obtained by means of a procedure discussed in Fama and MacBeth (1973): Cross-Sectional Factors.

Consider the cross-section regressions:

$$R_{it} = R_{zt} + R_{MCt}MC_{it-1} + R_{BMt}BM_{it-1} + R_{OPt}OP_{it-1} + R_{INVt}INV_{it-1} + e_{it} \quad (1.5)$$

Where the return of the  $i$ th security (portfolio) is regressed on a set of lagged characteristics. Here

- $R_{it}$  is the return of asset  $i$  at time  $t$ .
- $R_{zt}$  is the common return of all assets not captured by the characteristics

This specification exhibits a peculiar property: it can be proven that if each characteristic on the right-hand side ( $RHS$ ) is a Z-Score, the corresponding coefficient can be interpreted as the return on a zero investment portfolio for the  $LHS$  assets that sets the  $t - 1$  value of that characteristic to 1 and 0 out each other characteristic. Moreover  $R_{zt}$  equals the return of a equally weighted portfolio for the  $LHS$  assets. To put things in perspective, given the simplified specification

$$R_{it} = \alpha + R_{BMt}BM_{it-1} + \epsilon_i$$

if the  $BM$  characteristic is standardized in the sense that each observation is a deviation from the cross-sectional mean of that characteristics divided by its cross-sectional standard deviation,  $R_{BMt}$  can be seen as the return of a portfolio that goes long on assets with high value of the corresponding characteristics and short on those with low value, where the weights for each assets are the (standardized) characteristics' scores. A detailed proof can be found in appendix.

A key distinction between time-series (TS) and cross-sectional (CS) factors lies in the CS method's ability to incorporate other observable characteristics directly into the model equation, thereby isolating a "pure" factor effect. As discussed extensively in Lioui and Tarelli (2022), TS factors may unintentionally capture exposures to additional characteristics beyond the primary one of interest. This issue arises because TS factors often rely on sorting techniques that lack the capacity to fully separate these overlapping exposures, due to characteristics' correlations.

Conversely, the CS approach mitigates these issues by enabling the inclusion of multiple characteristics within its regression framework. By adjusting for other characteristics (e.g., size, book-to-market, profitability) in the CS regression, the resulting factor is "filtered", ensuring that the estimated factor returns are driven solely by the specific characteristic under examination.

When regression (1.5) is estimated and stacked across  $t$ , it can be viewed as an asset pricing model. In this perspective, by moving  $R_{zt}$  to the  $LHS$ , and swapping loadings with factors we obtain:

$$R_{it} - R_{zt} = MC_{it-1}R_{MCt} + BM_{it-1}R_{BMt} + OP_{it-1}R_{OPt} + INV_{it-1}R_{INVt} + e_{it} \quad (1.6)$$

With this specification, we have an asset pricing model where *LHS* returns are in excess of  $R_{zt}$ , the factors are the optimized returns estimated for each  $t$ , while the loadings are the pre-specified time-varying observed characteristics. Fama and French emphasize the superiority of this model over those considered in Fama and French (2020). The reason why this model shows better performance compared to the others is twofold: firstly, the presence of time-varying loadings is arguably a more natural setting for this type of models, furthermore, the characteristics result of adopting an optimization procedure with respect to the observed characteristics for each point in time, therefore the better performance in contrast to the Fama and French factors that are returns on pre-specified portfolios, and are not themselves optimized.

The factors obtained from (1.5) can also be used as an alternative set for model (1.4). In this setting, we now have a five-factor model with returns in excess of the risk-free rate and optimized factor apart from the market return:

$$R_{it} - R_{ft} = \alpha_i + \beta_{1i}(R_{mt} - R_{ft}) + \beta_{2i}R_{MCt} + \beta_{3i}R_{BMt} + \beta_{4i}R_{OPt} + \beta_{5i}R_{INVt} + \epsilon_{it} \quad (1.7)$$

This is the model (or a sub-specification of it) which the empirical results of *Chapter 6* are based on. The summary statistics from which the superiority of this model over model (1.4) is assessed can be found in the next chapter.

## 1.4 Two Pass Procedure

The *Two Pass Procedure* presented in Fama and MacBeth (1973) allows for the testing of asset pricing models like those mentioned above. As the name suggests it is a two-step technique whose main focus is to estimate *risk premia* associated with factors.

This procedure served as the inspiration for the cross-sectional factor methodology discussed earlier.

In the first step, given  $N$  assets and  $K$  distinct factors, *risk exposures* ( $\beta_{ik}$ ) corresponding to asset  $i$  and factor  $k$  are estimated running TS regression for each asset:

$$R_{it} - R_{ft} = \alpha_i + \beta_{i1}F_{1t} + \beta_{i2}F_{2t} + \dots + \beta_{iK}F_{Kt} + \epsilon_{it} \quad (1.8)$$

These regressions are performed via OLS obtaining a set of  $N$  estimates  $\hat{\beta}_{ik}$  for each of the  $K$  factors.

The second step consists of a cross-section (CS) regression where the averages with respect to  $T$  of the (excess) returns of each asset  $\overline{R}_i$  are regressed on the previously estimated  $\hat{\beta}_{ik}$ :

$$\overline{R}_i = \gamma_0 + \gamma_1 \hat{\beta}_{i1} + \gamma_2 \hat{\beta}_{i2} + \dots + \gamma_K \hat{\beta}_{iK} + \eta_i \quad (1.9)$$

thereby obtaining  $K$  distinct estimates of the *risk premia* ( $\gamma_k$ ) associated with each of the factors.

## 1.5 Factor Strength

The inference (and the identification) relative to the risk premia estimates that are retrieved through the two-pass procedure is heavily influenced by a variety of elements. One element of primary interest is the *strength* of the factors included in the model. In particular the inference becomes unreliable when the factor of interest is deemed to be weak.

Bailey et al. (2021) propose an estimator for factor strength based on its degree of "pervasiveness" identified by the number of its non-zero loadings, or more precisely, by the rate at which the number of non-zero loadings increases relative to the total number of loadings  $n$ . A factor is said to be "strong" (equivalently, strength equal to 1) if all its associated loadings are nonzero, the factor is said to be "semi-strong" if the rate lies between  $\frac{1}{2}$  and 1, while the factor is considered to be weak if the rate is less than  $\frac{1}{2}$ .

Considering a framework similar to what discussed in the previous sections:

$$R_{it} - R_{ft} = c_i + \beta_i F_t + \epsilon_{it} \quad (1.10)$$

where  $c_i$  is the unit specific effect,  $\epsilon_{it}$  is an idiosyncratic error and  $\beta_i$  is the factor loading for unit  $i$ .

The strength of factor  $F$  is given by the rate at which  $\omega_n^2 = \sum_{i=1}^n \beta_i^2$  rises with  $n$ , denoting the latter with  $\alpha$ . It is important to note that without further restriction on the factor loadings it is not possible to correctly identify  $\alpha$  when the factor is weak, but in most empirical applications the values of  $\alpha$  associated to the factors of interest are greater than  $\frac{1}{2}$  nonetheless.

To illustrate the estimation strategy, we consider the single factor model (1.10) where  $T$  observations are given on  $n$  cross-sectional units. The factor loadings are considered

to be nonzero for the first  $\lfloor n^\alpha \rfloor$  units and zero for the remaining units (where  $\lfloor \cdot \rfloor$  denotes the integer part function). Thus for some  $c > 0$ :

$$\begin{aligned} |\beta_i| &> c \quad \text{a.s. for } i = 1, 2, \dots, \lfloor n^\alpha \rfloor, \\ |\beta_i| &= 0 \quad \text{a.s. for } i = \lfloor n^\alpha \rfloor + 1, \lfloor n^\alpha \rfloor + 2, \dots, n, \end{aligned} \quad (1.11)$$

where  $\alpha$  is the strength of the factor  $F_t$ . The first step to estimate  $\alpha$  is to run  $n$  (OLS) TS regressions for each of the  $n$  units obtaining:

$$R_{it} - R_{ft} = c_{iT} + \hat{\beta}_{iT} F_t + \hat{\epsilon}_{it} \quad (1.12)$$

Let  $t_{iT} = \frac{\hat{\beta}_{iT}}{s.e.(\hat{\beta}_{iT})}$  be the  $t$ -statistic corresponding to  $\beta_i$  and  $\hat{\pi}_{nT} = n^{-1} \sum_{i=1}^n \hat{d}_{i,nT}$  the proportion of the regressions with statistically significant coefficients, where  $\hat{d}_{i,nT} = \mathbb{I}[|t_{iT}| > c_p(n)]$ .

The critical value function is given by:

$$c_p(n) = \Phi^{-1}\left(1 - \frac{p}{2n^\delta}\right) \quad (1.13)$$

In this context  $p$  is the nominal size of the individual tests,  $\delta$  is a parameter that rules the dependence of the critical value with respect to  $n$  (greater  $\delta$  leads to more conservative tests) and  $\Phi(\cdot)$  is the cumulative distribution function of the normal standard distribution. The authors suggest that, for the cases of interest where the true value  $\alpha_0 > \frac{1}{2}$ , a reasonable range for  $\delta$  is  $[\frac{1}{4}, \frac{1}{2}]$  or it can be chosen according to any cross-validation procedure.

The proposed estimator of factor strength is:

$$\hat{\alpha} = \begin{cases} 1 + \frac{\ln \hat{\pi}_{nT}}{\ln n}, & \text{if } \hat{\pi}_{nT} > 0, \\ 0, & \text{if } \hat{\pi}_{nT} = 0. \end{cases} \quad (1.14)$$

Bailey et al. (2021) justify the use of the exponent  $\alpha$  as a measure of factor pervasiveness, arguing that  $\alpha$  provides a more robust and discriminative indicator compared to a simple proportion  $\pi$  of significant loadings. While  $\pi$  quantifies the actual number of significant loadings, it tends to decay asymptotically with increasing  $n$ , making it challenging to interpret in large-dimensional settings. In contrast,  $\alpha$  reflects the growth rate of non-zero loadings, remaining stable as  $n$  grows and allowing for a more precise assessment of factor ‘‘strength.’’

Importantly,  $\alpha$  is *super-consistent* for factors with  $\alpha > 0.5$ , ensuring rapid convergence to the true pervasiveness of the factor and thus enhancing the reliability of

inferences regarding its identification.

To test  $H_0 : \alpha = \alpha_0$  the following score statistics is used:

$$z_{\hat{\alpha}:\alpha_0} = \frac{(\ln n)(\hat{\alpha} - \alpha_0) - p(n - n^{\hat{\alpha}})n^{-\delta-\hat{\alpha}}}{\left[p(n - n^{\hat{\alpha}})n^{-\delta-2\hat{\alpha}}\left(1 - \frac{p}{n^{\delta}}\right)\right]^{1/2}} \quad (1.15)$$

The null hypothesis is rejected if  $|z_{\alpha}| > cv$ , where  $cv$  is the critical value of the standard normal distribution at the desired significance level. The finite sample properties of the proposed estimator have been investigated using a number of Monte Carlo simulations. Figure 1.1 and 1.2 illustrate some of the results obtained by Bailey et al. (2021).

$n \setminus T$	Bias ( $\times 100$ )					RMSE ( $\times 100$ )					Size ( $\times 100$ )				
	60	120	200	500	1000	60	120	200	500	1000	60	120	200	500	1000
$\alpha_{10} = 0.75, \alpha_{20} = 0.85$															
100	1.13	1.18	1.15	1.07	1.03	1.65	1.54	1.52	1.43	1.40	9.00	4.10	3.65	2.40	2.25
200	1.46	1.46	1.39	1.32	1.32	1.68	1.62	1.55	1.47	1.47	14.50	9.50	8.30	7.10	6.60
500	1.28	1.30	1.21	1.15	1.13	1.41	1.37	1.28	1.22	1.20	22.40	13.55	10.00	8.05	8.20
1000	1.27	1.25	1.20	1.12	1.10	1.36	1.30	1.24	1.16	1.14	26.30	15.00	11.45	7.70	6.40
$\alpha_{10} = 0.80, \alpha_{20} = 0.85$															
100	0.63	0.67	0.65	0.61	0.58	1.13	1.00	1.00	0.95	0.92	27.60	18.10	18.55	17.95	19.70
200	0.90	0.97	0.91	0.89	0.86	1.11	1.10	1.05	1.01	0.98	20.45	12.60	11.45	9.75	7.75
500	0.82	0.90	0.85	0.82	0.80	0.94	0.96	0.90	0.87	0.86	22.00	12.45	8.30	7.00	7.10
1000	0.78	0.85	0.81	0.76	0.75	0.87	0.88	0.84	0.79	0.77	28.35	17.05	11.60	8.40	6.65

FIGURE 1.1: Bias, root mean square error (RMSE) and size ( $\times 100$ ) of estimating different strengths of the first factor in the case of two observed factors and non-Gaussian errors, when the strength of the second factor is set to 0.85. Source: (Bailey et al., 2021)

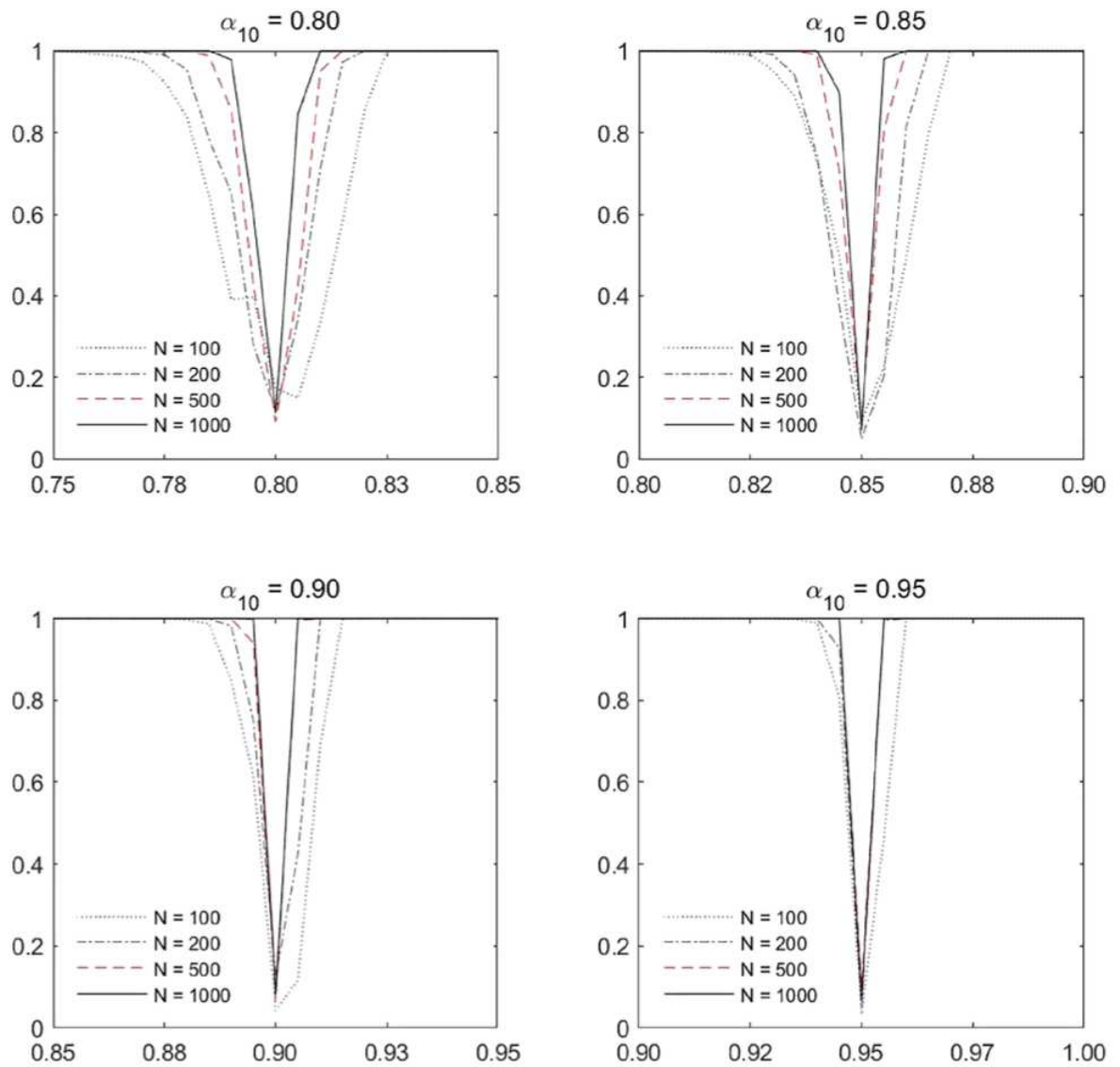


FIGURE 1.2: Empirical power functions associated with testing different strengths of the first factor in the case of two observed factors and non-Gaussian errors, when the strength of the second factor is set to 0.85,  $n = 100, 200, 500, 1000$ , and  $T = 200$ . Power is computed under  $H_1 : \alpha_{1a} = \alpha_{10} + k$ , where  $k = -0.05, -0.045, \dots, 0.045, 0.05$ . The number of replications is set to  $R = 2000$ . Source: (Bailey et al., 2021)



# Chapter 2

## Literature Review (B)

### 2.1 Empirical Results of the Fama-French Five-Factor Model

This section provides an overview of the empirical findings and explores the performance of each factor in the Fama-French five-factor model, with a particular emphasis on the Small Minus Big (SMB) factor.

Fama and French (2015) test the five-factor model across a range of portfolios formed on firm characteristics such as size, book-to-market ratio (B/M), profitability, and investment. The model's primary objective is to capture patterns in average returns more effectively than the three-factor model. The inclusion of the profitability (RMW) and investment (CMA) factors is intended to account for variations in returns that the size and value factors alone cannot explain. The results demonstrate that the five-factor model generally reduces unexplained return variation across portfolios when compared to the three-factor model, especially for portfolios with strong tilts toward profitability and investment.

#### 2.1.1 Focus on the SMB Factor

The size factor, or SMB (Small Minus Big), measures the return differential between small-cap and large-cap firms, capturing the size effect in equity returns. This factor is constructed by averaging returns on small-cap portfolios and subtracting the average return on large-cap portfolios, based on size and other sorts.

In Fama and French's analysis, the SMB factor remains robust across different sorting schemes, including 2x3 and 2x2 setups for size and B/M, as well as profitability

and investment sorts. Empirical tests show that the average monthly return for SMB is around 0.29% to 0.30% across the various constructions, with standard deviations between 2.87% and 3.13%. The factor’s strong correlations across different sorting methods indicate its stability and robustness as a size effect proxy.

However, the addition of the RMW and CMA factors alters the role of SMB in explaining size-related anomalies. With profitability and investment included, the five-factor model sometimes reduces the explanatory power of SMB, suggesting that part of the size effect may have been capturing variations in profitability or investment. Despite this, SMB remains a significant factor, with returns more than two standard errors from zero, reinforcing its contribution to the model.

When applied to portfolios with specific size and investment or profitability characteristics, SMB’s performance highlights the complexities in interpreting the size effect. In particular, Fama and French find that SMB displays stronger explanatory power for portfolios tilted toward small-cap stocks with high profitability or low investment, consistent with the size effect’s role in capturing risk-return dynamics among small firms. In the context of the five-factor model, SMB aligns closely with traditional interpretations, but its interaction with RMW and CMA requires careful consideration in analyzing overall model performance.

### 2.1.2 Summary

In summary, the empirical results affirm the robustness of the five-factor model in describing average returns across various portfolio sorts. The SMB factor, while still relevant, interacts with RMW and CMA in ways that suggest a nuanced understanding of the size effect when controlling for profitability and investment. Future research may further investigate these interactions, particularly in markets where size-based strategies are popular.

<b>Factor</b>	<b>Mean Return (%)</b>	<b>Std Dev (%)</b>	<b>t-Statistic</b>
Market (RM-RF)	0.50	4.49	2.74
Size (SMB)	0.29	3.07	2.31
Value (HML)	0.37	2.88	3.20
Profitability (RMW)	0.25	2.14	2.92
Investment (CMA)	0.33	2.01	4.07

TABLE 2.1: Summary statistics for the monthly factor returns (1963-2013).

## 2.2 Empirical evidence for Time-Series Models using Cross-Sectional Factors (Fama and French, 2020)

### 2.2.1 The framework

Fama and French (2020) compare the performance of asset pricing models that utilize factors estimated from cross-section regressions (cross-section or CS factors) versus those estimated using the traditional time-series approach (time-series or TS factors). Specifically, as discussed in Chapter 1, they employ the Fama and MacBeth (1973) cross-sectional regression methodology to construct CS factors that are then used in a time-series asset pricing framework to evaluate their performance compared to standard TS factors.

The empirical analysis is based on monthly returns from July 1963 to August 2018, covering 662 monthly returns. The cross-section factors are constructed from the 18 value-weighted (VW) portfolios used to produce the traditional TS factors, including portfolios formed on characteristics like size (MC), book-to-market (BM), profitability (OP), and investment (INV). In this way, Fama and French ensure a level playing field for comparing CS and TS factors by using the same underlying portfolios.

Two models are particularly considered: Model (1.7), which uses CS factors, and Model (1.4), the Five-Factor time-series model by Fama and French (2015)

### 2.2.2 Comparisons and summary statistics

The main empirical finding is that the model using cross-section factors outperforms the time-series factor model across a range of test portfolios. The cross-section model better captures average returns on a diverse set of portfolios formed on characteristics like market capitalization, book-to-market ratio, profitability, and investment.

The use of CS factors yields lower pricing errors when explaining average returns. The analysis in the paper shows that the average absolute pricing error ( $A|a|$ ) for Model (1.6) is consistently lower compared to the pricing error for Model (1.3).

The CS factors exhibit better performance due to the fact that they are optimized month by month based on the cross-section of asset characteristics. This leads to more accurate factor loadings and a reduction in pricing errors.

Table 2.2 presents the summary statistics for the selected models. These results pertain to portfolios formed based on characteristics such as size, book-to-market ratio, profitability, and investment. The table provides several metrics to evaluate the model performance:

- $Sh^2(a)$ : Represents the maximum squared Sharpe ratio for the intercepts. Here,  $a$  is the vector of intercepts obtained from the model, and  $\Sigma$  is the covariance matrix of the regression residuals. Calculated as  $Sh^2(a) = a'\Sigma^{-1}a$ , this metric assesses the efficiency of the intercepts in relation to their residual covariance.
- $A|a|$ : The average of the absolute values of the intercepts for the left-hand side (LHS) assets. This metric provides a simple measure of the model's pricing errors by taking the absolute intercept values without considering covariances.
- $A|t(a)|$ : The average of the absolute values of the  $t$ -statistics for the intercepts. This statistic measures the significance of the intercepts across the cross-section, offering insight into the pricing error magnitudes.
- $Aa^2/V\bar{r}$ : Represents the average of the squared intercepts for the LHS assets, divided by the cross-sectional variance of LHS average returns. This ratio estimates the proportion of cross-sectional return dispersion that the model fails to capture.
- $A\lambda^2/V\bar{r}$ : An adjusted version of  $Aa^2/V\bar{r}$  that accounts for noise in the intercept estimates. Here,  $\lambda^2 \equiv a^2 - s^2(a)$ , where  $s^2(a)$  is the squared standard error of each intercept. This metric adjusts the proportion of return dispersion missed by the model by removing the estimated noise in the intercepts.
- $AR^2$ : The average  $R^2$  across all regressions. This statistic provides a summary measure of the model's goodness of fit.
- $As(a)$ : The average of the standard errors of the intercepts across the cross-section. This measure helps assess the variability in the intercept estimates.
- $As(e)$ : The average of the standard deviations of the regression residuals, which provides insight into the model's overall residual variability.

These metrics collectively provide a comprehensive assessment of model performance, with Model (1.6) showing better performance compared to Model (1.3), indicating that the cross-section factors offer superior explanatory power in this context.

Model	$ \alpha $	$ t(\alpha) $	$A^2/V_f$	$\Delta\lambda^2/V_f$	$AR^2$	$As(\alpha)$	$As(\epsilon)$	$Sh^2(\alpha)$	GRS	$p(GRS)$
Panel A1: LHS assets are 75 without-momentum portfolios; factor loadings are constant regression slopes										
(3) $R_m - R_f, SMB, HML, RMW, CMA$	0.081	1.34	0.33	0.22	0.92	0.060	1.48	0.356	2.53	0.000
(4) $R_m - R_f, R_{Mc}, R_{BM}, R_{OP}, R_{INV}$	0.070	1.09	0.24	0.12	0.92	0.063	1.54	0.347	2.45	0.000
(3) $R_m - R_f, SMB, HML, RMW, CMA, UMD$	0.072	1.17	0.26	0.15	0.92	0.061	1.48	0.340	2.35	0.000
(4) $R_m - R_f, R_{Mc}, R_{BM}, R_{OP}, R_{INV}, R_{MOM}$	0.068	0.98	0.24	0.10	0.91	0.068	1.63	0.317	2.14	0.000
Panel A2: LHS assets are 100 without- and with-momentum portfolios; factor loadings are constant regression slopes										
(3) $R_m - R_f, SMB, HML, RMW, CMA, UMD$	0.082	1.29	0.21	0.14	0.92	0.063	1.52	0.594	2.95	0.000
(4) $R_m - R_f, R_{Mc}, R_{BM}, R_{OP}, R_{INV}, R_{MOM}$	0.083	1.15	0.22	0.13	0.91	0.069	1.65	0.596	2.89	0.000

TABLE 2.2: Summary of intercepts from regressions explaining returns in excess of  $R_f$ , models use CS Factors (Model 3), and Fama-French Factors (Model 4). Source: (Fama and French, 2020)

## 2.3 Factor Strength (Pesaran and Smith, 2021)

The paper by Pesaran and Smith (2021) investigates the implications of pricing errors and the varying strength of factors on the estimation of risk premia using the Fama-MacBeth two-pass estimator. The authors employ a comprehensive empirical framework to evaluate the performance of the five Fama-French factors, considering different levels of factor strength.

### 2.3.1 Data and Methodology

The empirical analysis uses the Fama-French Five-Factor model, which includes factors such as market, SMB (Small Minus Big), HML (High Minus Low), RMW (Robust Minus Weak), and CMA (Conservative Minus Aggressive). To address the shortcomings deriving from the instability of factor strength over time, the authors apply a rolling window approach of 10 years to estimate factor strength for each factor over time. The analysis covers the period from September 1989 to May 2018, resulting in 345 rolling estimates for each factor. The rolling window approach provides insights into the temporal variations in factor strength.

They use all stocks in the S&P 500 portfolio that have at least 10 years of return history, for each month from September 1989 to May 2018. The list is updated monthly and includes at least 400 stocks, with an average number of 442 stocks. This procedure avoids the possible survivorship bias caused by the changing composition of S&P 500 portfolio.

The study makes use of a modified two-pass estimation procedure, focusing on factor strength as a key determinant of model performance. Factors are classified as strong, semi-strong, or weak based on the proportion of securities with statistically significant factor loadings.

### 2.3.2 Results: Factor Strength Estimation

Table 2.3 presents the rolling estimates of factor strength for the five Fama-French factors. The results indicate that only the market factor consistently maintains a high level of strength, with an average value close to one, and can therefore be classified as *strong*. In contrast, the other four factors generally exhibit lower strengths, with average values below 0.8, namely *semi-strong* factors. This suggests that the explanatory power of the non-market factors may be limited, particularly for smaller sub-samples.

TABLE 2.3: Rolling Estimates of Factor Strength for Fama-French Factors: September 1989 - May 2018

	<b>Factor</b>	<b>Average</b>	<b>Minimum</b>	<b>Maximum</b>	<b>Standard Deviation</b>
Market	0.99	0.95	1.00	0.01	
SMB	0.72	0.60	0.80	0.05	
HML	0.68	0.50	0.76	0.07	
RMW	0.65	0.48	0.74	0.08	
CMA	0.67	0.52	0.73	0.06	

The table shows that the strength of the market factor is consistently strong, whereas the remaining factors have lower strength values, indicating that they influence fewer securities compared to the market factor. The temporal variation in factor strength, as depicted by the rolling estimates, highlights the dynamic nature of the factors' influence on asset returns.

Figure 2.1 illustrates the evolution of factor strength over time for each of the five Fama-French factors. The market factor remains strong throughout the sample period, whereas the other factors exhibit significant fluctuations. Figure 2.2 focuses on the SMB factor. The rolling estimates reveal that the strength of the non-market factors varies substantially, implying periods where these factors have a limited impact on asset pricing.

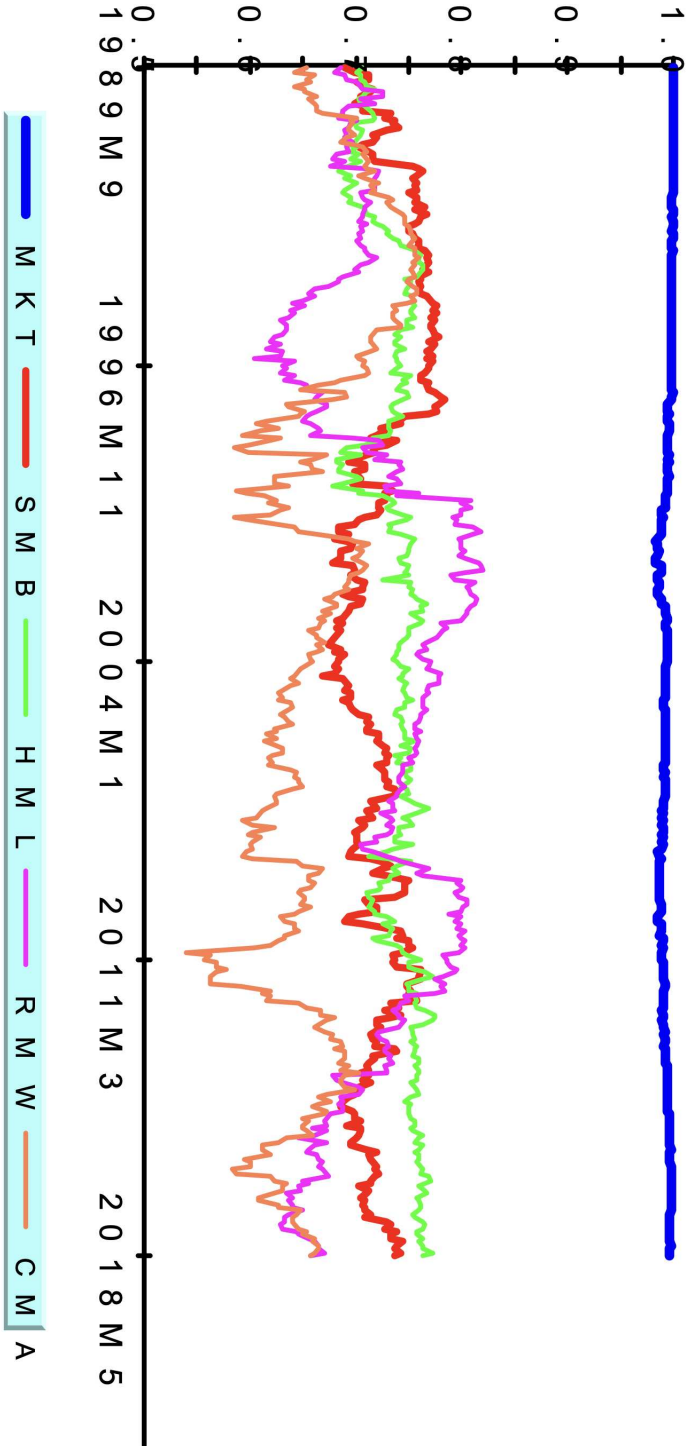


FIGURE 2.1: Rolling Estimates (10 years windows) of Factor Strength for the Five Fama-French Factors (September 1989 - May 2018). Source: (Pesaran and Smith, 2021)

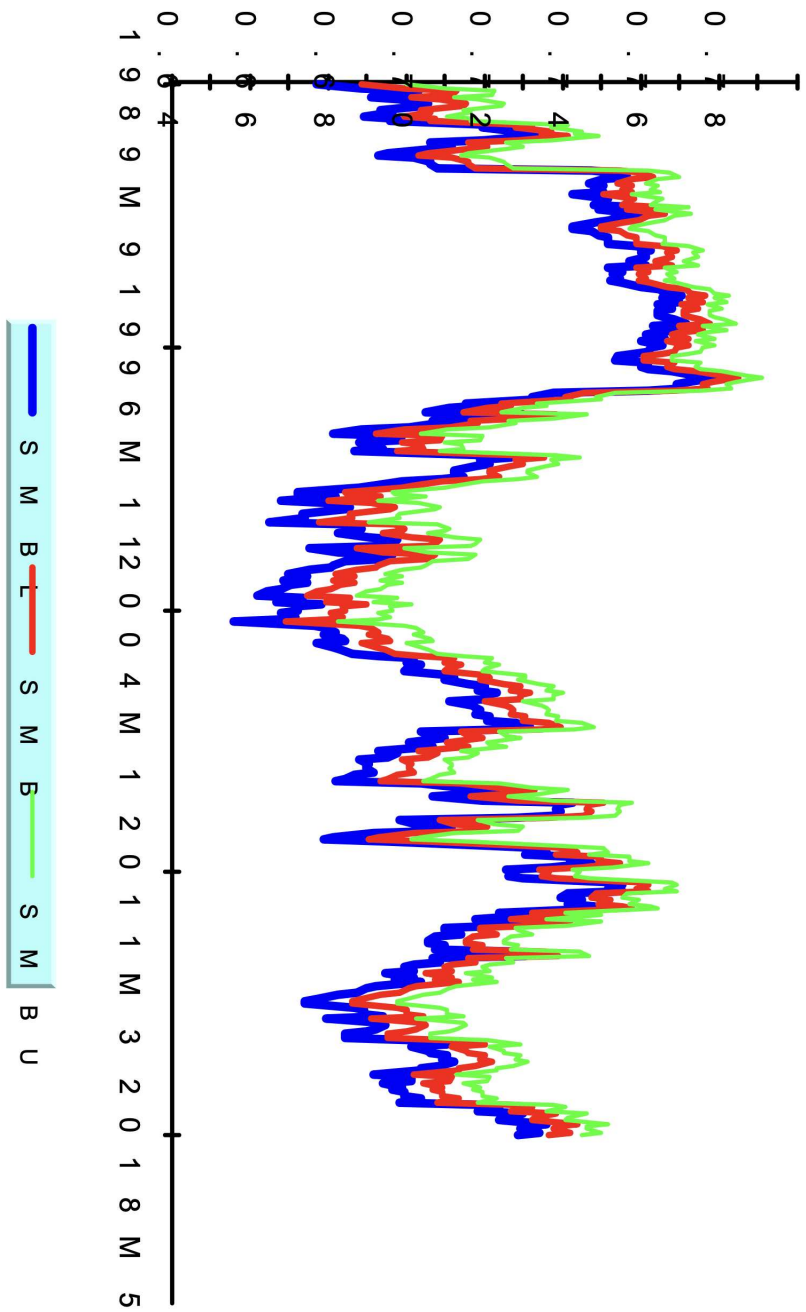


FIGURE 2.2: Rolling Estimates (10 years windows) of Factor Strength for SMB Factor (September 1989 - May 2018). Source: (Pesaran and Smith, 2021)



# Chapter 3

## Data

### 3.1 Data collection

All data adopted in the present work have been downloaded from LSEG Refinitiv, the data have been recovered from a collection of lists and data request tables. A global perspective has been adopted, downloading equity prices and company-specific variables from the financial markets of the following countries:

- **Europe:** Austria, Belgium, Denmark, Finland, France, Germany, Ireland, Italy, Netherlands, Norway, Poland, Portugal, Spain, Sweden, Switzerland, United Kingdom;
- **Africa:** South Africa;
- **America:** Brazil, Canada, Chile, Mexico, United States of America;
- **Asia:** China, India, South Korea, Taiwan.

For every country, the list of companies quoted in the reference market (as identified by Refinitiv) has been downloaded, with the exclusion of the USA where both the New York Stock Exchange and the NASDAQ are considered. In selecting the companies, for which both equity data as well as other variables (including balance sheet data) are downloaded, the attention is restricted to companies whose equities are classified as *Major Securities* and *Primary Quote*. In addition, equities are quoted in the local currency for each market. In order to avoid the survivorship bias, both *Dead* and *Active* equities are selected. Finally, the data frequency is set to monthly (end of month values) from the 31st of December 1979 to the 31st of May 2024, giving a total of 534 monthly observations.

The analyses will be developed in a common currency, the United States Dollar, and therefore the exchange rates between the local currencies and the USD have also been recovered in the same time range adopted for the equities. When available over the required sample (31/12/1979-31/05/2024) the WM Reuters exchange rates is selected, otherwise other exchange rates available in Refinitiv are chosen.

Overall, for each company, the following variables are available:

- **P**: the equity price;
- **RI**: the total return index;
- **MV**: market value;
- **BV**: book value of equity;
- **OI**: operating income, used with **BV** to calculate the scaled operating profit ( $OI/BV$ );
- **TA**: total assets used to calculate investments as the one-period lagged logarithmic growth rate of **TA**;
- **PE**: price/earnings;
- **DY**: dividend yield;
- **MTBV**: market to book value;
- **EPS**: earnings per share;
- **EPS12**: earnings per share 12 months forward;
- **EPSSUR**: earnings per share mean surprise;
- **LEV**: leverage ratio;
- **TD**: total debt;
- **EBITDA**: earnings;
- **NETDEBT**: debt minus cash (net financial position), used as a ratio to **EBITDA** to calculate a leverage ratio;
- **NETINT**: net interests, used as a ratio to **EBITDA** to calculate a leverage ratio;

- Additional cross-sectional information: delisting date, country of incorporation (derived from the market), and the economic sector (using a classification in 14 sectors plus a residual class including companies which are unclassified).

The data have been downloaded in the form of Excel files (including macro, due to the adoption of Refinitiv *request table* tool), and before proceeding to the analyses, all data have been imported into R for a set of preliminary checks and data cleaning procedures.

## 3.2 Pre-Processing

The raw data consists of records from 66,059 companies as detailed in Tables 3.1 and 3.2. In a preliminary step, all the companies for which the company name was not reported in Refinitiv were excluded from the dataset; these cases are in general associated with companies without valid data for the selected time range. Following, as the focus on this work is on the small cap effect, which is associated to the so-called *Size* factor put forward by Fama and French (1993), based on the companies market value, all the cleaning procedures have been first based on the availability of the market value data as well as on the availability of the equity price data (both in the form of equity prices as well as in terms of total return indexes). Focusing on the market value, and given the presence of both dead and active equities over a pre-specified time range (31/12/1979-31/05/2024), it is possible that some equities do not have valid data (i.e. equities whose trading activity was ended before 12/1979). These invalid data might be represented both as missing values as well as a flat market value of the considered time range, thus all equities with fully missing market value or with a constant market value from December 1979 to May 2024 are removed.

At the company level, data availability is also associated with the trading of companies in the financial market. For companies classified as *Dead* Refinitiv frequently provides a flat market value and/or a flat price from the month where market trading of a given equity has stopped (the delisting date). Therefore, at the end of each series, data characterized by flat values in the market value and/or in the equity price, including the delisting month are removed.

It is worth mentioning that data from a few companies appeared in multiple geographical areas. To address this issue, the duplicated records of these companies have been retained exclusively in the market where the company is legally headquartered.

The companies are originally classified into the following economic sectors: Energy, Basic Materials, Consumer Staples, Financials, Industrials, Utilities, Real Estate, Health

Care, Telecommunications, Consumer Discretion, Technology, and Unclassified. The sectors Unclassified, Unquoted Equities, Other Equities, and Suspended Equities have been merged into a new “Unclassified” sector, which incorporates all companies that were initially assigned to any of these categories.

Finally, due to the need of converting all variables into a common currency, some additional checks are performed. First, in a given financial market, companies reporting data in a different currency are excluded, as they most likely correspond to companies based in a different country whose equities are traded in the considered financial market. Second, for countries belonging to the Euro-area, given the downloaded time range covers a period where local currencies were used, it is essential to verify that all values are reported in Euro. If that is not the case, Euro-based time series are recovered adopting the fixed Euro conversion rates. After performing the common currency conversion, data from Brazil prior to January 1999 were excluded to address certain inconsistencies resulting from the currency conversion.

The literature of factor estimation by means of cross-sectional approaches usually includes additional cleaning rules. The most diffused one is the exclusion of the so-called *penny stocks*, companies whose equity price, in a given cross-section, is below 5 USD (this threshold is a common choice). Another common choice is the use of winsorization, that is, the exclusion of companies falling in the lower  $\alpha\%$  or in the upper  $1 - \alpha\%$  of the empirical distribution for at least one of the variables included in the analysis. We are using a tailored combination of these two approaches. First of all, we exclude penny stocks. Second, we do not use in a straight way a windorisation approach. In fact, as one of the purposes of the present work is to dig into the role played by huge companies (in particular of the *Maginificent Seven*) (MAG7), winsorization at the market value level would certainly exclude the MAG7, that would be located in the upper extreme tail of the market value distribution. Therefore, for the Market Value, we opted for the exclusion of a small number of outliers graphically identified by looking at the data and at preliminary evaluations of the quantities of interest. Clearly, this is an ad hoc procedure, but would allow maintaining full control of the excluded companies and of their number, at least in some evaluations. This winsorization scheme lead us to exclude, from the very beginning, 5 companies (one for US and 4 located in European markets). Then, for all other variables, we adopt a fixed winsorization scheme, dropping the 0.5% in each tail (i.e., overall, we discard in each cross-section, 1% of companies).

After the exclusion of penny stocks and these 5 companies, the dataset cross-sectional

dimension on a monthly basis evolves as reported in Figure 3.5.<sup>1</sup> In addition, in Figure 3.6, we report the percentage of penny stocks that have been excluded from the original dataset at each point in time. Their number is increasing over time, growing with the data coverage and the number of overall companies included in the dataset. In relative terms, this suggest that the fraction of companies traded in the financial markets with a small equity price has increased over the years (once all prices have been converted in US Dollars). We note that, in the following, graphs of the dataset coverage will be included whenever additional variables will be taken into account for the estimation of the Size factor. A robustness check on the impact of excluding penny stocks with different thresholds (thus deviating from the 5 USD threshold commonly adopted in the academic literature) is included in the following section.

### 3.3 Data coverage

This section describes the financial markets covered in the following analyses, the corresponding data coverage as recovered from the data provider adopted, and the data coverage that will be used in the following sections.

At the end of the cleaning procedure, the database composition by country and by economic sector is reported in Tables 3.1 and (3.2). Notably, the full cleaned dataset includes more than 58 thousands companies for which data on market value and equity price are available at least for part of the sample. At the country level, as expected, the more represented countries are the United States, the United Kingdom, India and China. The Euro area countries, when combined (Finland, France, Germany, Ireland, Austria, Italy, Netherlands, portugal, Belgium, Spain), have a number of companies larger than the United Kingdom.

At the sector level, a large fraction of companies (however, sensibly reduced when focusing on the cleaned data) is not accompanied by the sector information in LSEG Refinitiv. These companies are, in general, dead companies. The most represented sectors are the Industrial, the Consumer Discretionary, the Financial and the Technology.

Figure 3.1 report the database coverage for two reference variables, the equity price and the market value, and for each month in the sample. As the database include dead companies, the full coverage (100%) is never attained. The two plots show that the coverage tends to increase over time, due to the increasing number of companies quoted in financial markets. However, the coverage at the beginning of the sample is quite low,

---

<sup>1</sup>The figure scale has been adopted to allow comparison of the cross-sectional dimension with other cases covered in the present report where, by considering additional factor drivers or control variables, the cross-sectional dimension might sensibly reduce.

	Raw		Cleaned	
	N.	%	N.	%
Brazil	829	1.25%	642	1.10%
Canada	3502	5.30%	3229	5.54%
Chile	474	0.72%	335	0.57%
China	5734	8.68%	5376	9.22%
Denmark	614	0.93%	432	0.74%
Finland	408	0.62%	336	0.58%
France	2595	3.93%	2371	4.07%
Germany	1818	2.75%	1537	2.64%
India	8119	12.29%	6789	11.65%
Ireland & Austria	466	0.71%	407	0.70%
Italy	971	1.47%	898	1.54%
Korea	4108	6.22%	4043	6.94%
Mexico	412	0.62%	317	0.54%
Netherlands	515	0.78%	435	0.75%
Norway	868	1.31%	780	1.34%
Poland	1288	1.95%	1248	2.14%
Portugal & Belgium	821	1.24%	545	0.94%
South Africa	1146	1.73%	1020	1.75%
Spain	636	0.96%	531	0.91%
Sweden	1693	2.56%	1549	2.66%
Switzerland	574	0.87%	473	0.81%
Taiwan	2950	4.47%	2818	4.83%
United Kingdom	7524	11.39%	6606	11.33%
USA: NASDAQ	8036	12.16%	7822	13.42%
USA: NYSE	9958	15.07%	7749	13.29%
TOTAL	66059	100.00%	58288	100.00%

TABLE 3.1: Absolute and relative data coverage at the country level. In a couple of cases countries are combined for efficiency in the data download step.

being around 10%. Figures 3.2 and 3.3 show the coverage by country (in some cases countries are combined to avoid graphing patterns based on a small number of equities). The figures show that for some countries almost no data are available in the first ten year of the sample (up to the end of 1989). For market value, the absence of data is particularly severe for Brazil, China, Chile, Finland, India, Mexico, Poland, South Korea, Spain, and Taiwan. Limited data are also available for Belgium (included with Portugal), France, Germany, Italy, Norway, South Africa, Sweden, and for the United States if we consider only the NASDAQ market. Even if this last list of markets include data, the limited coverage indicates that only a small number of companies is available, thus leading to potential uncertainty in the cross-sectional analyses during the 80's.

A further issue, again impacting on data availability, is the conversion to a common

	Raw		Cleaned	
	N.	%	N.	%
Energy	2513	3.80%	2379	4.08%
Basic materials	5043	7.63%	4878	8.37%
Consumer Staples	3027	4.58%	2912	5.00%
Financials	9360	14.17%	8564	14.69%
Industrials	11164	16.90%	10718	18.39%
Utilities	1315	1.99%	1247	2.14%
Real Estate	2852	4.32%	2734	4.69%
Healthcare	4827	7.31%	4714	8.09%
Telecommunications	1445	2.19%	1408	2.42%
Consumer Discretionary	10033	15.19%	9647	16.55%
Technology	6877	10.41%	6769	11.61%
Unclassified	7603	11.51%	2318	3.98%
Total	66059	100.00%	58288	100.00%

TABLE 3.2: Absolute and relative data coverage at the sector level.

currency (not included in the previous plots). Figure 3.4 reports the availability of exchange rates, needed for the conversion into USD. Again, some exchanges rates are not available for the full sample, in particular those for Brazil, China, India, and Mexico, that would reduce the number of available companies in the cross-section for the first 15 years of the sample.

Combining these two elements, a proper choice would be to exclude from the analysis the 80's when focusing on the entire dataset. Differently, when considering only the US, data from the 80's will be used.

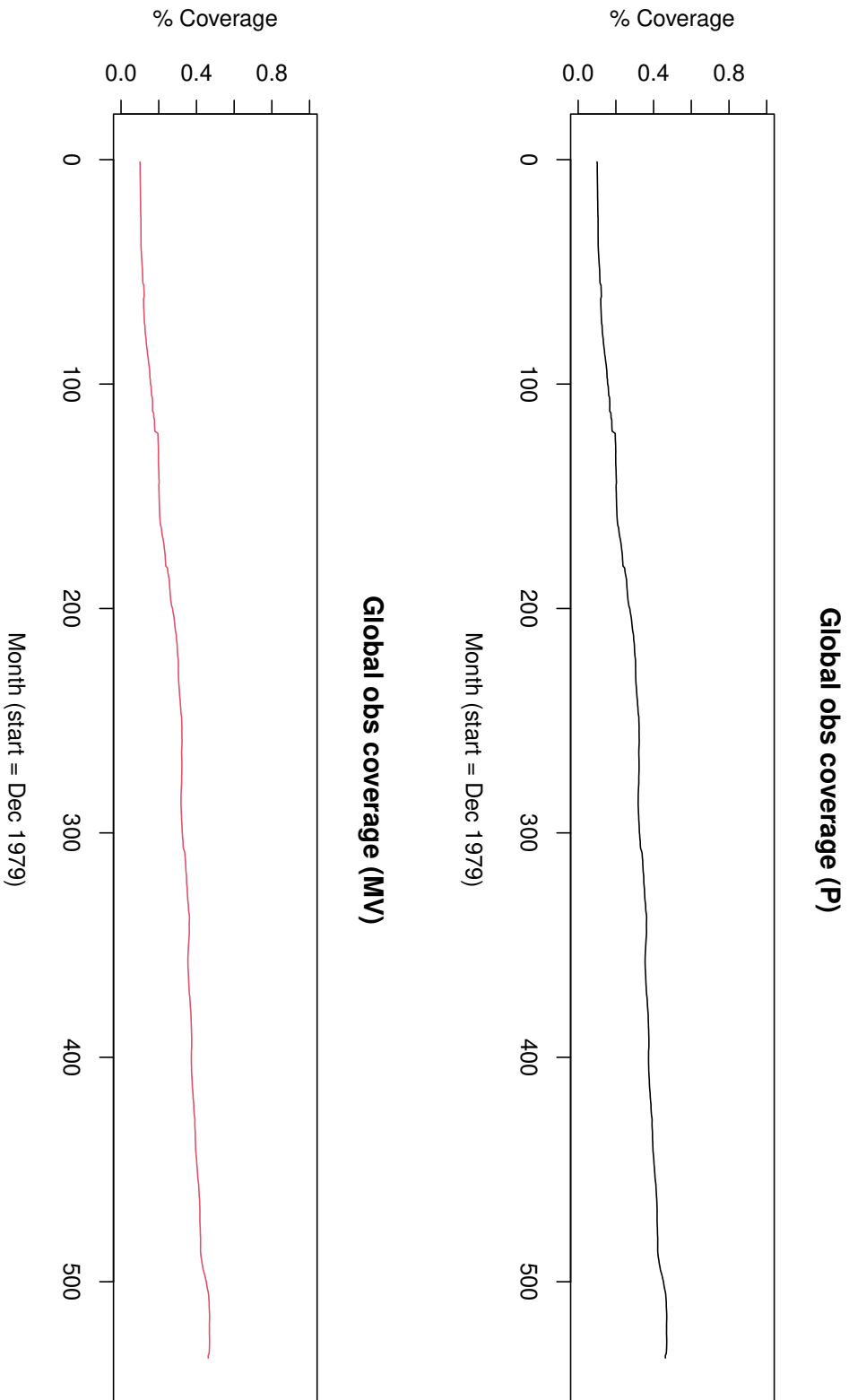


FIGURE 3.1: Percentage of available data in the global cross-section for each month and for both Price (upper panel) and Market Value (lower panel). The database include dead companies so total coverage is never attained.

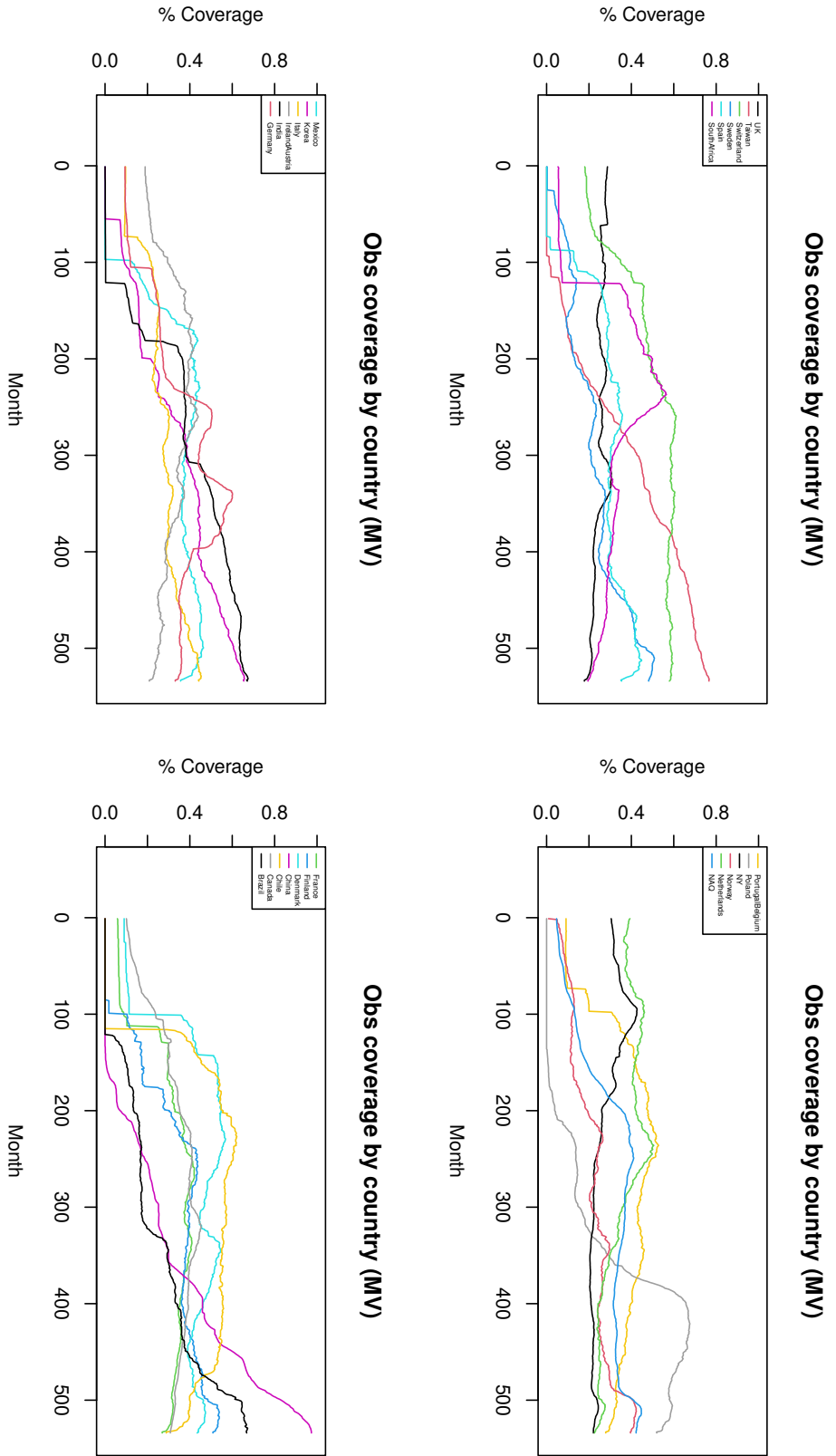


FIGURE 3.2: Percentage of available data at the country level for Market value. The database include dead companies so total coverage is never attained. Some countries are combined in a single line to avoid reporting patterns based on a small number of equities.

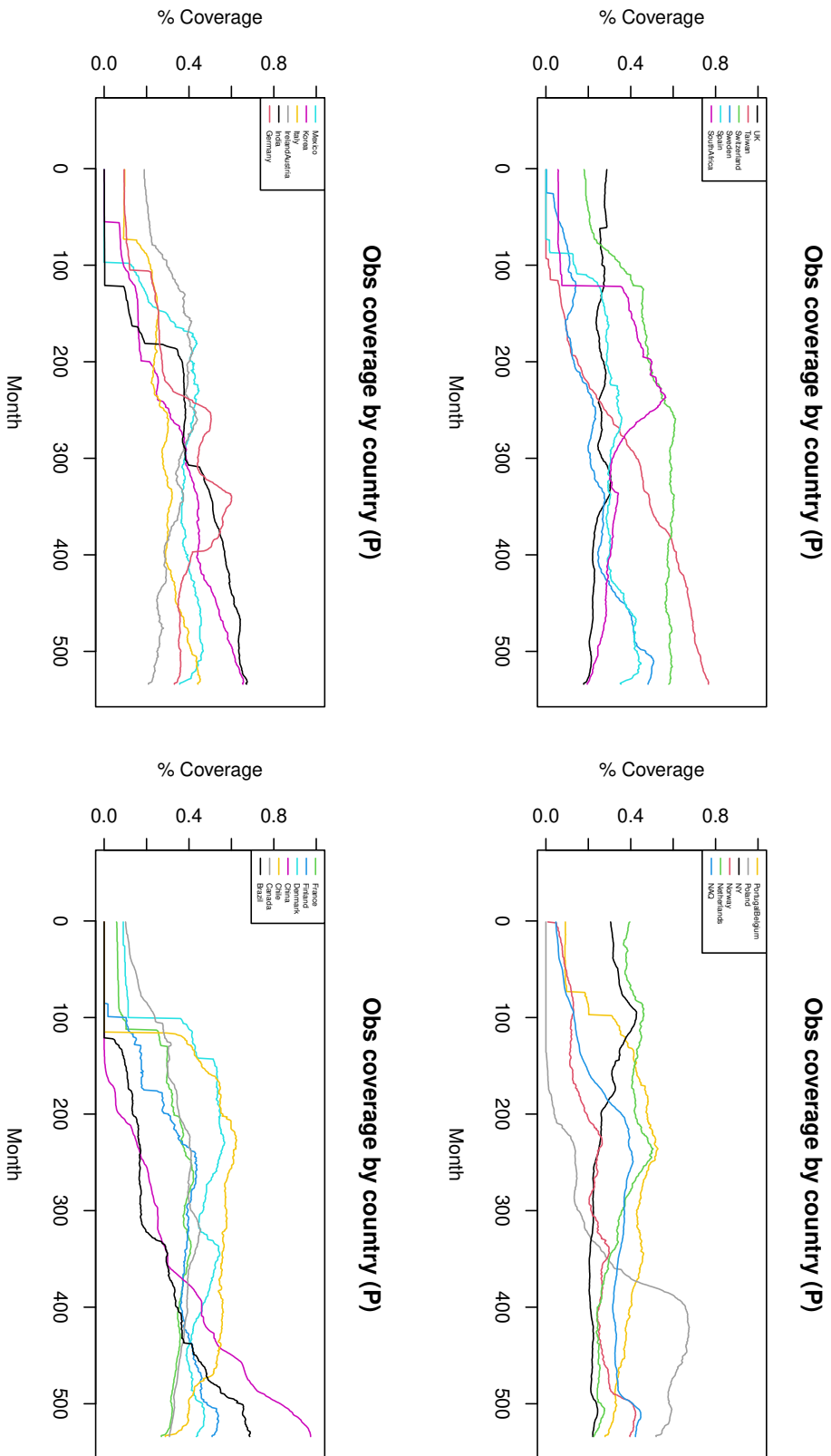


FIGURE 3.3: Percentage of available data at the country level for equity Price. The database include dead companies so total coverage is never attained. Some countries are combined in a single line to avoid reporting patterns based on a small number of equities.

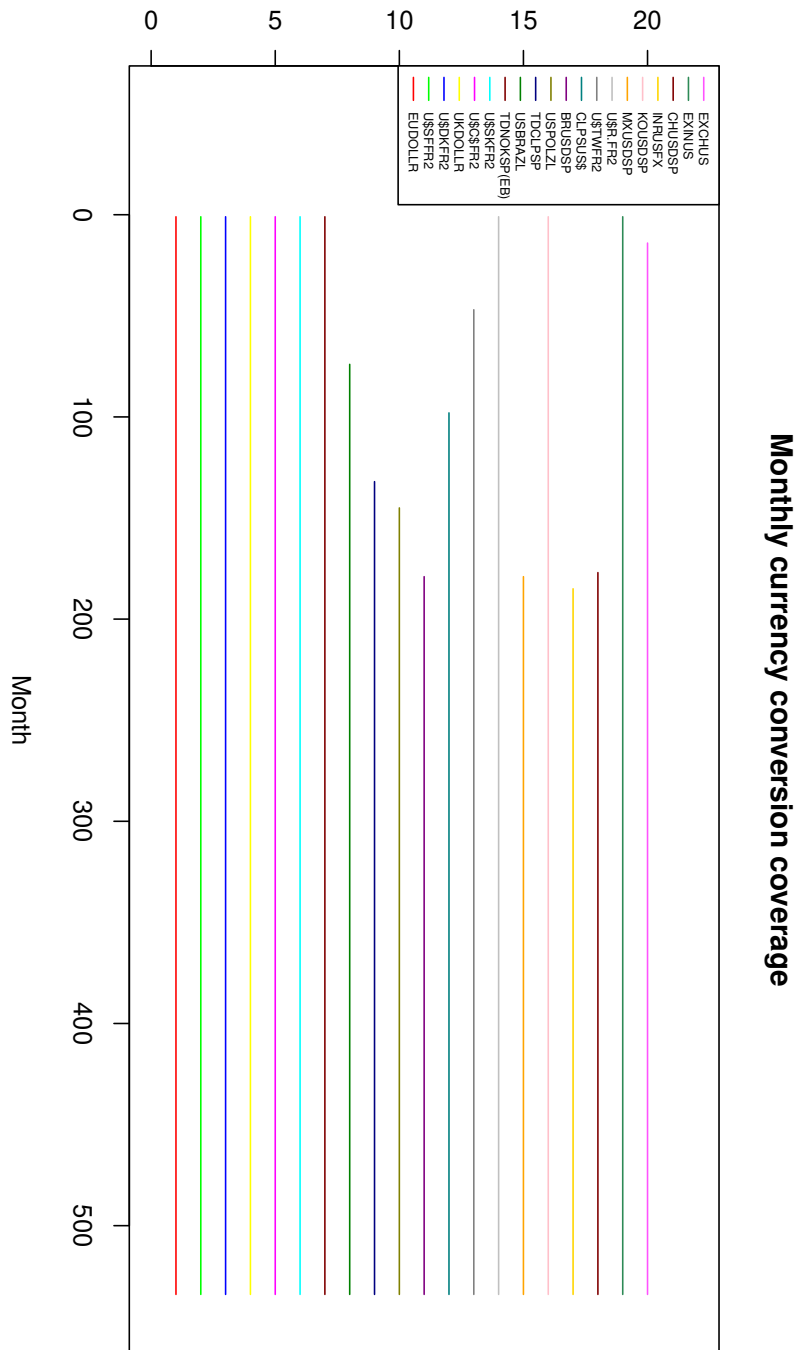


FIGURE 3.4: Availability of exchange rates with respect to the United States Dollar. The bars are in the same order as in the legend.

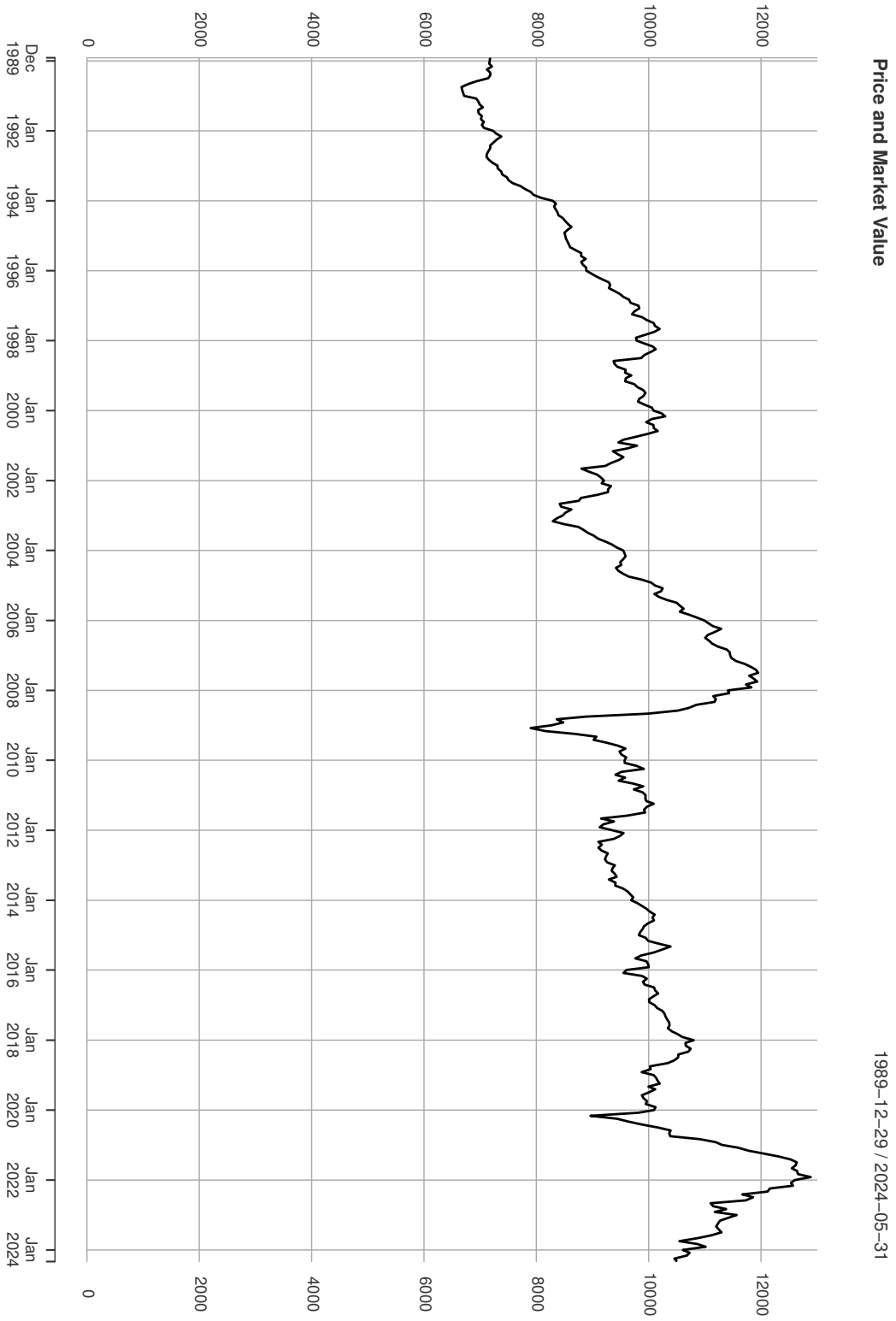


FIGURE 3.5: Dataset cross-sectional dimension after excluding penny stocks and considering companies with available Price and Market Value.

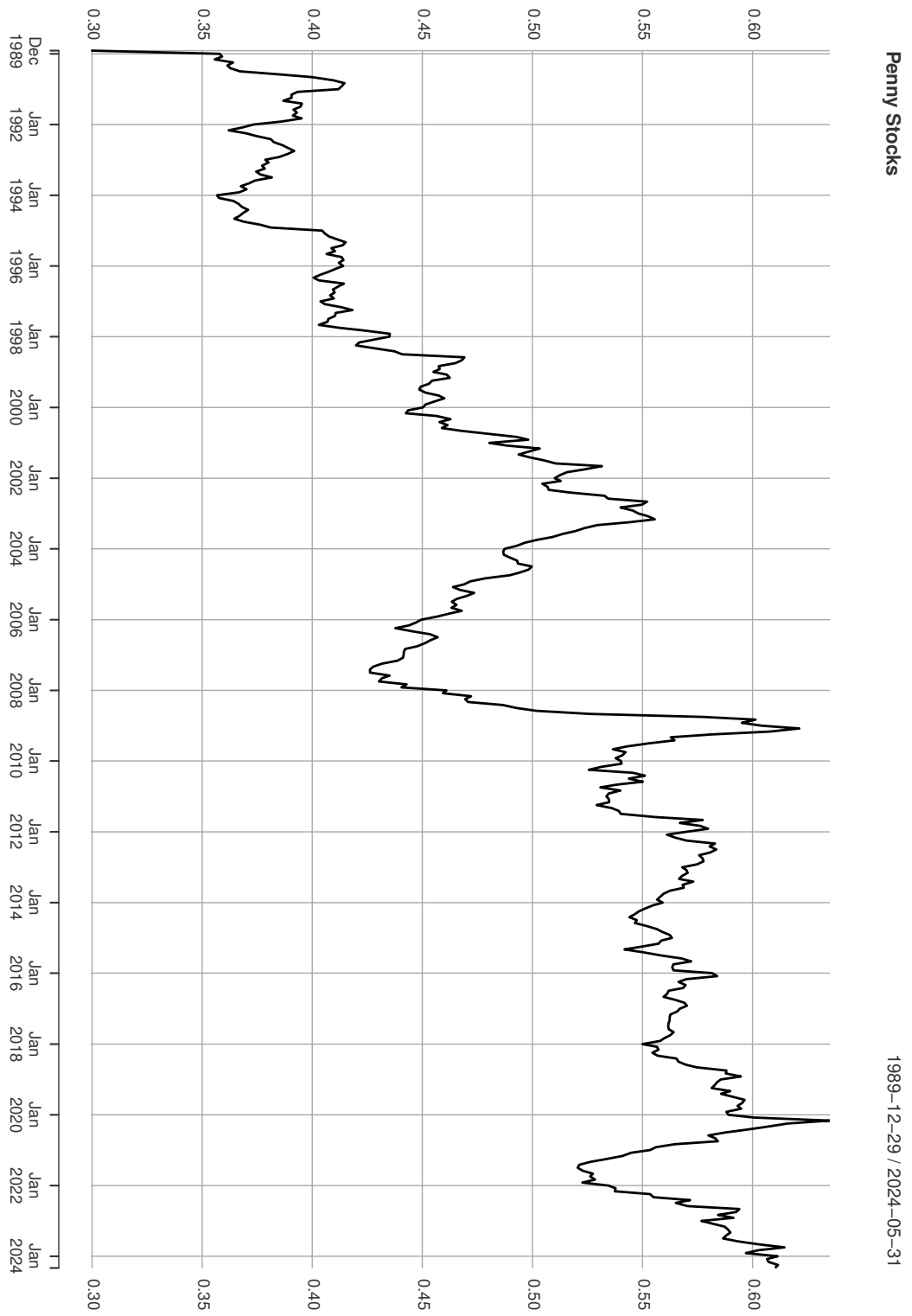


FIGURE 3.6: Percentage of penny stocks that have been excluded from the analysis at each cross-section.



# Chapter 4

## Cross-Sectional factors

### 4.1 The framework

This chapter presents the results of the estimation procedure used to retrieve the factors for the subsequent analyses.

Equation (1.5) can be written in more general terms as:

$$R_{it} = R_{zt} + \sum_{k=1}^K R_{kt} F_{kit-1} \sum_{c=1}^C \beta_{cit} X_{cit-1} + e_{it} \quad (4.1)$$

Where  $F_{kit}$  is the  $k$ th observed characteristics at time  $t$  for unit  $i$ ,  $X_{cit}$  is the  $c$ th variable for unit  $i$  from the set  $\mathbf{X}$  of control variables.

Given a model specification (choice of which characteristics to be included in the model), equation (4.1) fitted for a specific time  $t$  via ordinary least squares (OLS) yields factor estimates (or multiple factor estimates in the case where more than one characteristics is included). When the regressions are performed for each  $t = 1, \dots, T$ , the result is a time series (estimate) of monthly risk premia (i.e., the extra-performance provided by small cap stocks compared to large cap ones when considering the size characteristic), which will be used in a time series analysis in subsequent chapters.

It is important to note that, from this point onward, whenever a characteristic is mentioned, we are actually referring to its cross-sectionally (CS) standardized value. This allows for the interpretation of the factors as long-short portfolios, going long on assets with high value of the characteristic and short those which exhibit low value of the characteristic. The primary focus of this work is on the Size factor, which is associated with the Market Value characteristic, but due to the properties of OLS estimates, It is essential to verify the robustness of the factor estimates when additional characteristics

are included in the model, as discussed in detail in Chapter 6.

We expect that the size factor is at most weakly correlated with other factor (e.g. Fama & French factors) that may potentially be included in the model. This implies that the size factor estimates obtained from a model including multiple factors should not differ significantly from those obtained from a regression using only the Market Value (MV) characteristic.

To check the robustness of the *Size Factor* multiple regressions are performed with different specifications.

Several estimates resulting from different model specifications are showed in the following sections. For each specification, two primary settings are analyzed. The first involves applying the estimation procedure to different geographical sub-samples of the dataset (i.e., Europe, US, and Emerging Markets, including Brazil, India, China, and South Africa). This analysis covers the time window from 1990 to 2024, depending on data availability for the countries considered. The second setting involves applying the same procedure exclusively to the US sub-sample, covering the period from 1980 to 2024, to take advantage of the greater data availability characterizing US companies. It is important to note that when selecting a specific sub-sample, the standardization of characteristics is carried out by excluding all securities that are not part of the sub-sample of interest.

The signs of the *Market Value (MV)* and *Market-to-Book Value (MTBV)* characteristics are reversed when included in the regressions. The rationale behind this choice is that reversing the signs allows for a coherent comparison between the estimated factors and the Fama-French time-series (TS) factors. For example, consider the SMB factor, which, as discussed earlier, represents the extra return of small companies relative to large companies. Taking the negative of *MV (-MV)* results in an estimate that reflects the return of a long-short portfolio, where the portfolio goes long on companies with low MV values and short on companies with high MV values. In the following plots, we will mainly provide the cumulative value originated by a unit investment at the start of the period considered.

## 4.2 The baseline estimate

The first step is the estimation of the Size factor in the baseline case, that is, when this is the only estimated factor and there are no control variables. To perform this

estimation, the market value of equities for each company  $MV_{it}$  and the equity (total) returns are used,  $R_{it} = \frac{(RI_{it}) - (RI_{it-1})}{(RI_{it})}$ . Therefore, we have  $F_{it} = \{MV_{it}\}$  and  $X_{i,t} = \{\emptyset\}$ .

Results are reported in figures from 4.1 to 4.6. Figure 4.1 plots, for the entire dataset, the estimated cumulative monthly size premium (lower panel) from 1990. The plot shows some interesting patterns. First of all, three different regimes appear: a first low growth phase during the 90's, ending with a limited decline; a second high growth phase for the following 10-12 years (up to 2012), characterized by local instabilities; a last period of limited growth, taking a stabilizing pattern from 2021. Such behavior challenges the relevance and the impact of the size factor, in particular during the most recent years, and calls for a deeper evaluation on the possible role exerted by companies by country/geographical area, sector, and according either to the impact of additional variables as well as for the inclusion in the analysis of the largest (huge) companies.

If we consider the geographical slicing of the dataset, Figure 4.2, we note some differences. We separate the United States, the European countries, and the remaining markets which we cluster into an *Emerging* group (named BICS as it includes Brazil, India, China and South Africa). The pattern for the US closely follow that of the full (cross-sectional) sample, even if we acknowledge that the stabilizing pattern seems to start somewhat earlier (even from 2014), and we note a sudden increase in 2021, which we might attribute to the effects of the Covid-19 period. When moving to the European area, the high growth period continues until 2020, followed by a sudden increase as observed for the US, but then a declining pattern starts, thus suggesting a decrease in the relevance of the factor. An even different pattern appears for the emerging group, where we might still identify three phases, but with different growth rates: a first period, from 1990 to 2010 is characterized by a high growth in the size factors; a more stable period is observed from 2010 to 2021, with limited growth (in absolute terms), and then the factor growth restarts in the most recent years. The heterogeneity in the size factor growth among geographical areas might be due to a number of possible reasons: economic fundamentals (the economic growth rates in the three areas is known to be different); differentials in equity market risk; differences in the classification of companies when focusing on the full sample versus a *local* sample; an impact of the methodology followed for the factor construction given that the standardization of the Market Value is made within each slice of the dataset as oppose to the global standardization adopted before.<sup>1</sup> These element could explain both the different patterns as well as the different scales of the cumulative factor values. In fact, while the global factor reaches a total

---

<sup>1</sup>We note that the use of global standardization and a local factor construction is not possible as the parameters estimated by means of cross-sectional regressions will not correspond to the returns of a long-short portfolio.

Range	Full	US	Europe	BICS
1990-1994	4.43%	15.93%	1.93%	65.96%
1994-1999	2.00%	5.53%	0.25%	60.37%
2000-2004	13.64%	32.20%	8.21%	48.22%
2005-2009	6.73%	11.77%	3.02%	45.78%
2010-2014	3.90%	10.36%	3.43%	6.52%
2015-2019	2.07%	2.85%	3.30%	9.41%
2019-2024	1.38%	2.45%	3.98%	15.26%
1990-2024	38.91%	110.20%	26.56%	672.43%

TABLE 4.1: Cumulative returns by period and for the full sample - the range covers the period January 1990 to May 2024.

growth of about 40% in 35 years, the US growth is 110% for the US market, only 25% for the European markets and close to 700% for the emerging area. We stress that differences are large because the construction of the factors when slicing the dataset in different ways impact also on the standardization and the subsequent creation of long and short positions. As an example, a company based in an emerging market might be considered a small company when focusing on the entire dataset (thus ideally leading to a long position on that company) but it could become a large company when considering only emerging markets. Therefore, when constructing the size factor, the appropriate selection of the sample perimeter becomes fundamental to get a proper picture.

In Table 4.1 the cumulative returns by periods of 5 years and for the full sample are reported, distinguishing the total dataset from its decomposition by geographical area. We observe that the size factor is much less pronounced for the European countries compared to the other areas. In addition, both for the US and Europe, the factor returns are lower in the last ten years compared to what we have observed in the past.

Differentials due to sector characteristics emerge in a clear way when we consider a split of the dataset by economic sectors, Figure 4.3 and Figure 4.4. We remind that the factor is estimated restricting attention to a single sector. Therefore, the (cross-sectional) standardization of the Market Value made prior to the cross-sectional regression is made within the sector to maintain the interpretation of the estimated parameter as the size factor premium.

Finally, when focusing on the US market only, the larger data availability allows to extend backward the construction of the size factor up to 1980; Figure 4.5 reports the estimated factor. Notably, the overall pattern is confirmed and the additional ten years at the beginning belong to the slow growth period previously identified.

We close this section with a robustness check whose purpose is to identify the impact

Sector	Return	Sector	Return
Basic materials	193.40%	Consumer staples	96.03%
Consumer discretionary	38.68%	Energy	147.31%
Financials	53.51%	Health care	211.16%
Industrials	99.27%	Real estate	59.82%
Telecommunications	300.29%	Technology	35.14%
Utilities	217.25%	Unclassified	20.97%

TABLE 4.2: Cumulative total return of the monthly size factor premium by economic sector. The period ranges from January 1990 to May 2024.

of the penny stock threshold adopted for the evaluation of size factor. As previously explained, we use a 5 USD threshold coherently with previous studies to identify, and remove, penny stocks. In Figure 4.6 we plot the estimated Size factor obtained by varying the threshold from 1 to 5 USD with a 1 USD step. Notably, the only case in which the behavior of the factor sensibly differ, is at the 1 USD level, and only from 2001 onward. Differently, for all the other cases the factor is almost unaffected by the threshold. Given these evidence, we maintain the 5 USD threshold in the following analyses.

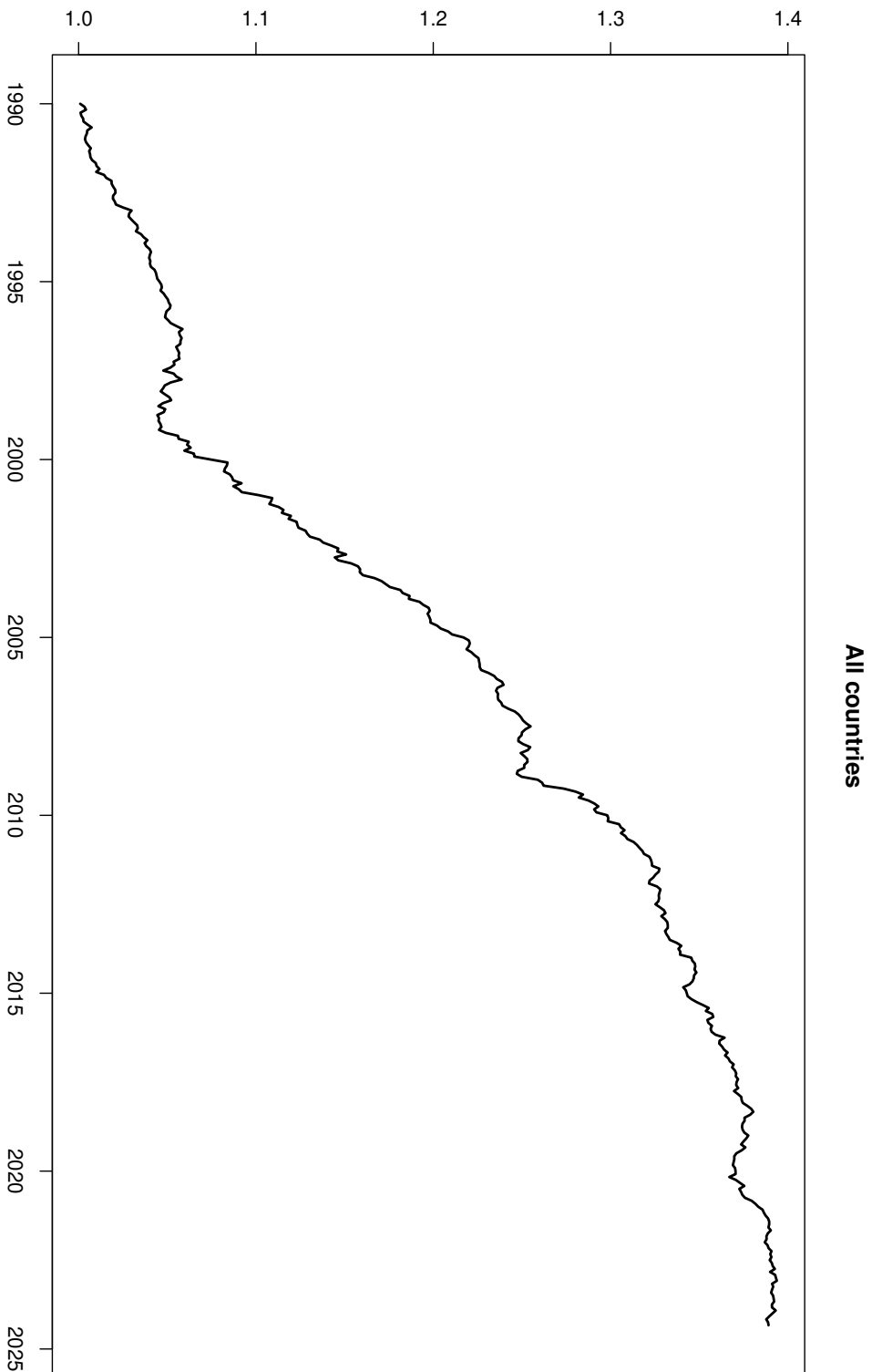


FIGURE 4.1: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as unique explanatory variable - data from 1990 - All countries.

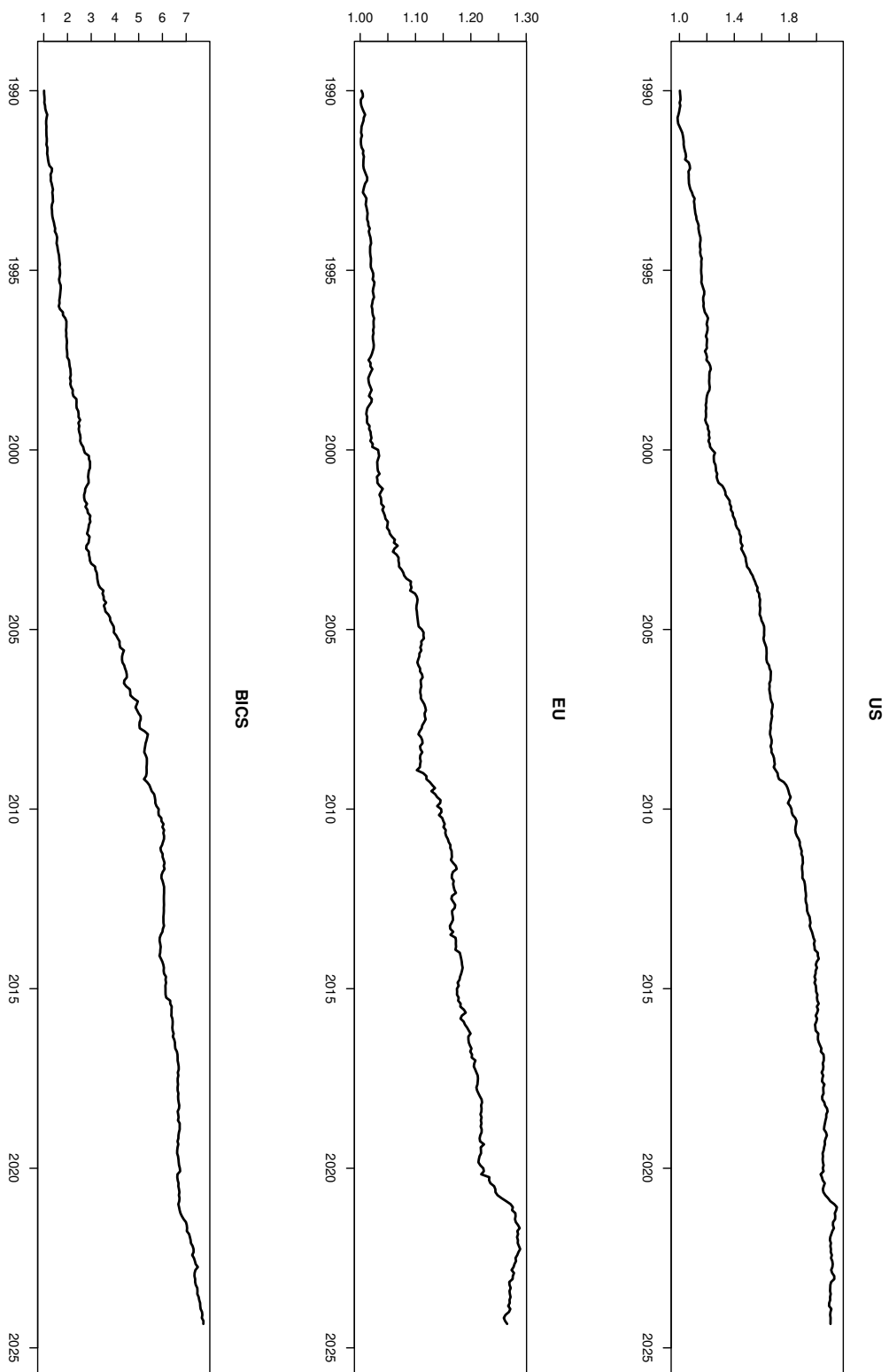


FIGURE 4.2: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as unique explanatory variable - data from 1990 - Estimates made at the geographical group level (i.e., the Size factor is estimated using companies belonging to a specific geographical area).

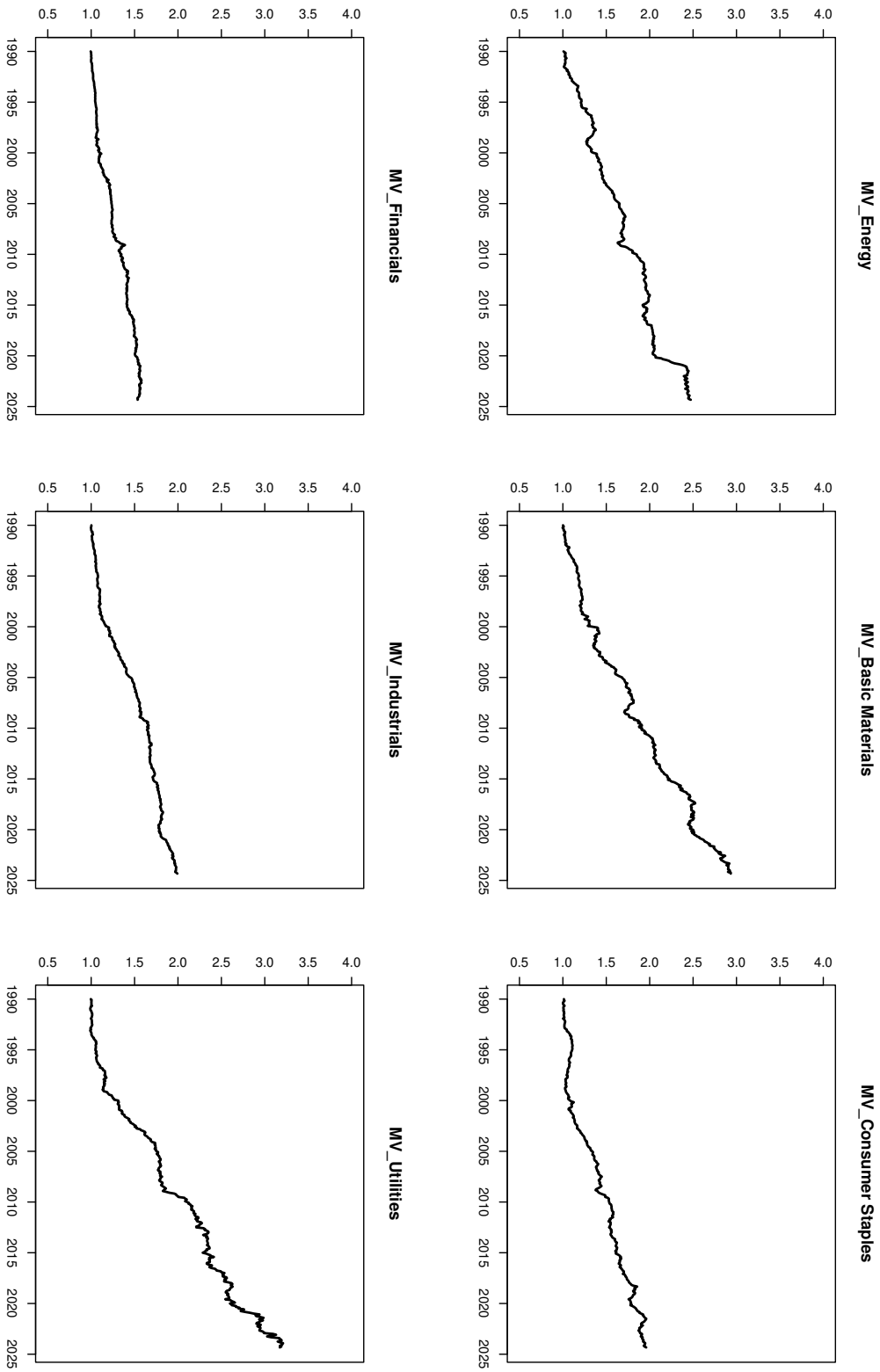


FIGURE 4.3: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as unique explanatory variable - data from 1990 - Estimates made at the sector group level (i.e., the Size factor is estimated using companies belonging to a specific sector) - Part A.

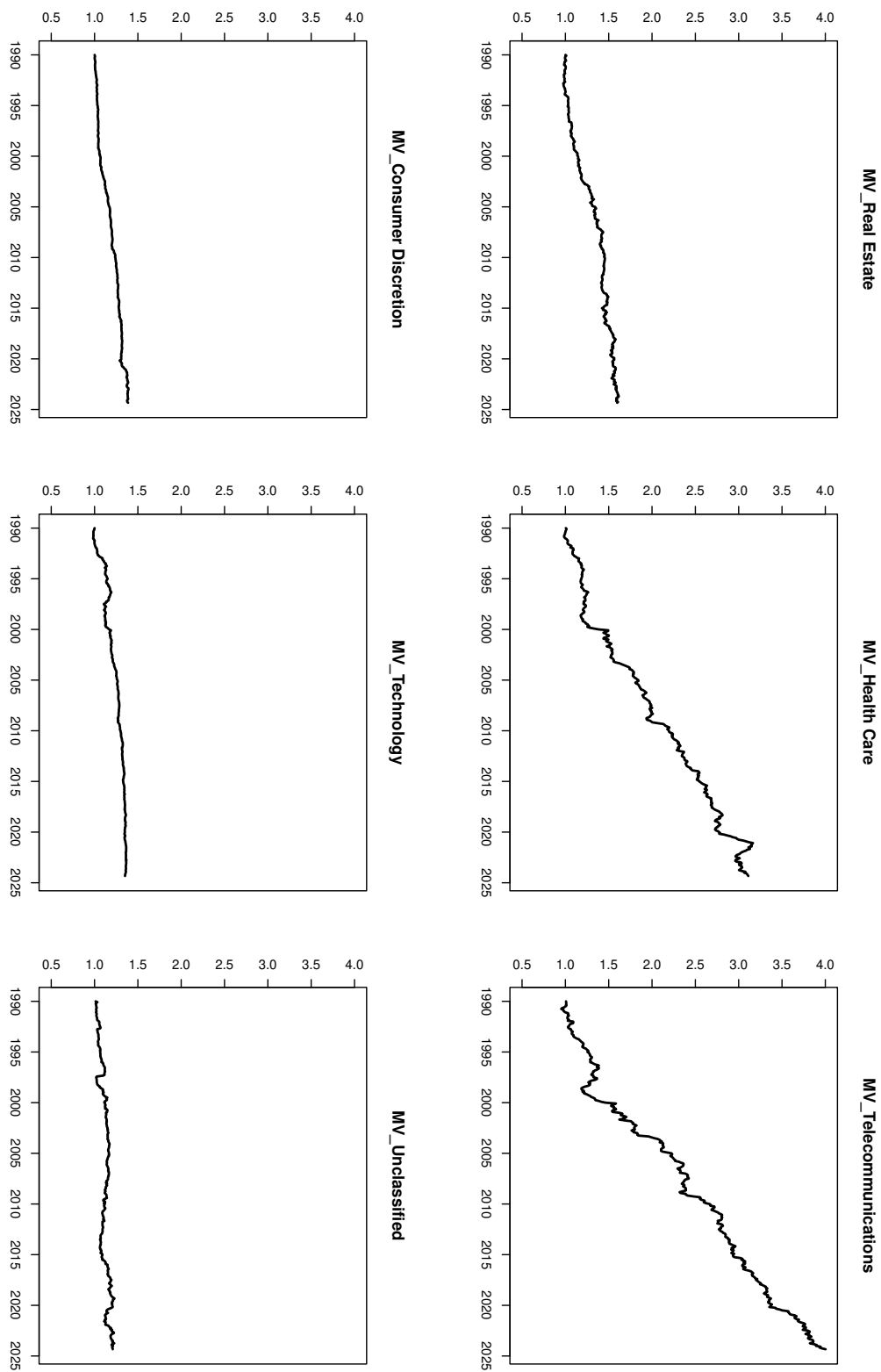


FIGURE 4.4: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as unique explanatory variable - data from 1990 - Estimates made at the sector group level (i.e., the Size factor is estimated using companies belonging to a specific sector) - Part B.

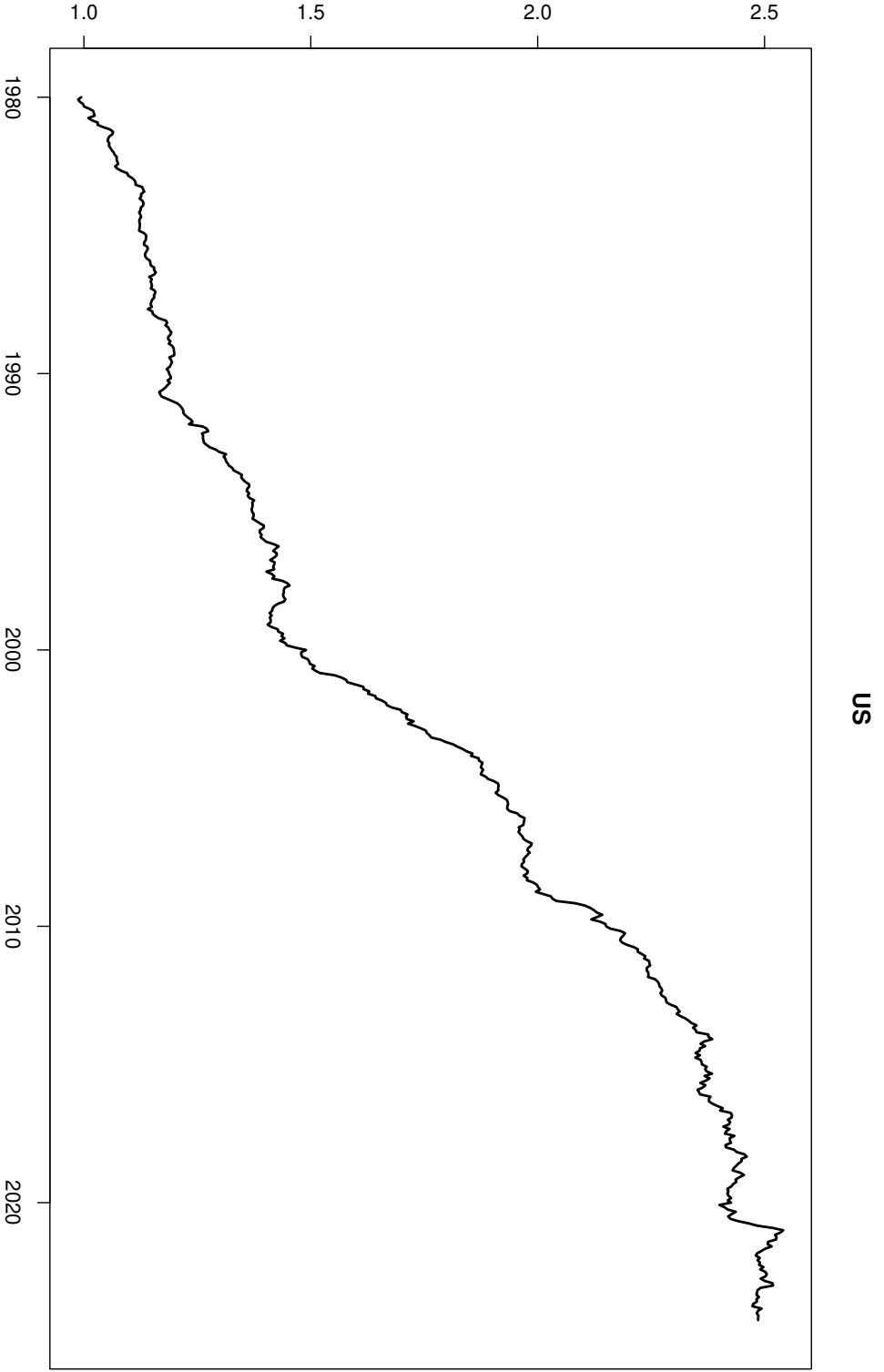


FIGURE 4.5: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as unique explanatory variable - data from 1980 - United States.

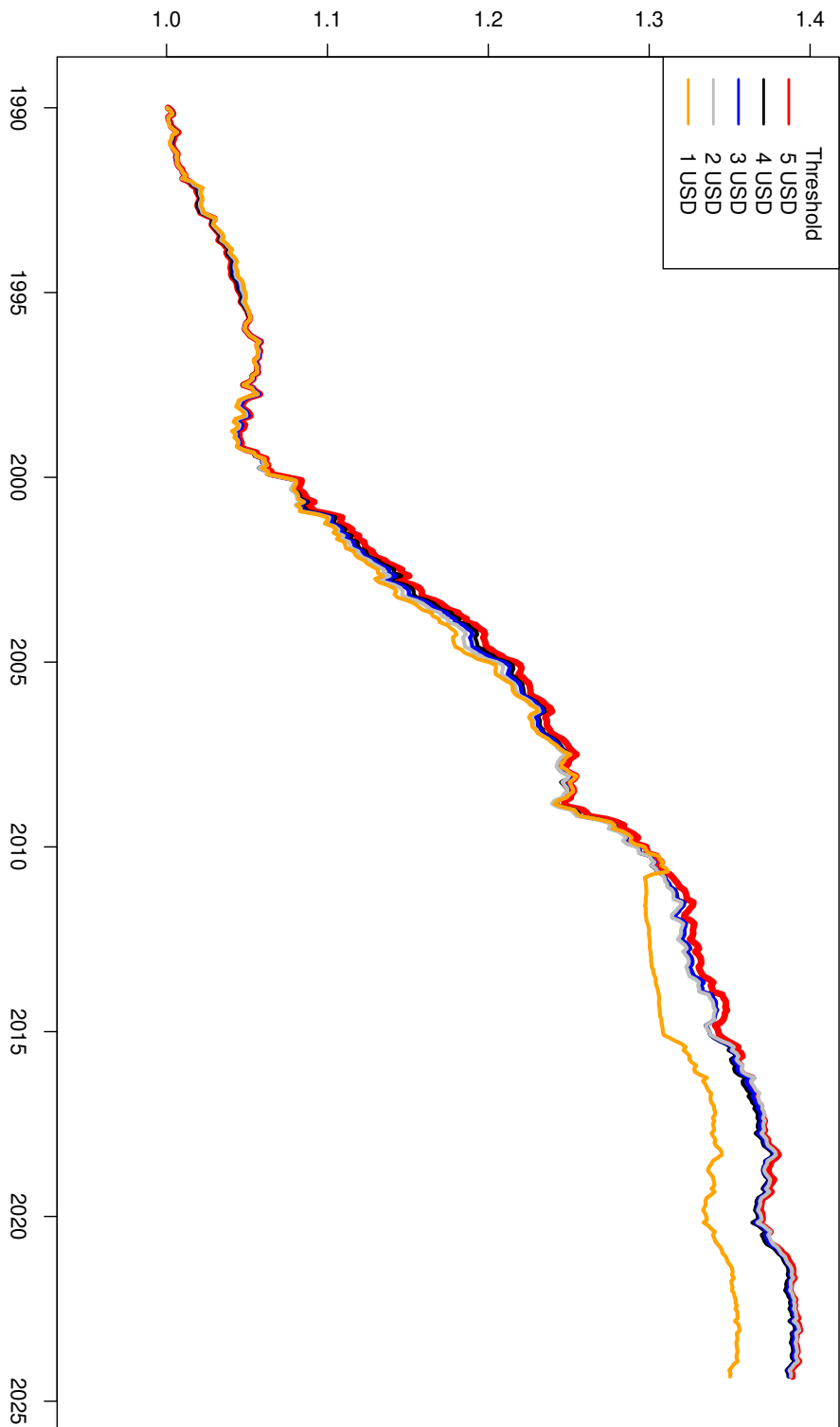


FIGURE 4.6: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as unique explanatory variable - data from 1990 - All countries - Monitoring the impact of different threshold for penny stocks selection.

## 4.3 Challenging the Small Cap Effect

### 4.3.1 Including other variables

Building on the previous analysis we now perform the estimation of the *Size* factor iteratively including multiple variables of interest into the model equations. The first variable included in the model is the MTBV characteristic, therefore jointly estimating the *Size* and the *Value* Factors. Note that these factors will be used in the following chapters as an alternative to the corresponding TS factors from Fama and French (2015) 3 Factor Model. Given that the authors were suggesting a factor mimicking portfolio going long on high Book-to-Market value and short on companies with low Book-to-Market value, as mentioned previously, we include the variable in the model with a negative sign in order to achieve coherence with its TS factor counterpart. In addition, we also make use of winsorization, cutting the top and bottom 0.5% of companies by MTBV (total cut, 1% of the companies at each cross-section) to mitigate the impact of extreme observations.

In this setting  $F_{it} = \{MV_{it}, MTBV_{it}\}$  and  $X_{it} = \{\emptyset\}$ .

Figure 4.7 plots the coverage of the dataset, considering those companies which have  $P$ ,  $MV$  and  $MTBV$  available simultaneously (those observations which have at least one missing variable at time  $t$  are omitted from the regressions). It is important to keep track of the coverage because adding a variable to the model specification could lead to a non-negligible change in the number of available observations.

The blue lines from Figures 4.10 to 4.14 show how the estimated size factor is affected by the introduction of the MTBV characteristic.

It is evident that the size factor, when estimated jointly with the value factor, exhibits a contraction in scale, while the general pattern remains largely unchanged. This is especially pronounced for the Emerging Markets compared to other geographical sub-samples, which exhibit a total return over the full sample that shrinks to roughly 300% (from 700% as seen previously). The only exception is the EU sub-sample, for which the estimated size factor appears to remain stable. It is also worth noting that, for the majority of the geographical subsets, the cumulative factor (and thus the cumulative return of the portfolio) at the beginning of the sample are comparable, with differences compounding over time.

Next, we included both the Investment Growth (IG) and Operating Profit (OP) characteristics into the model equations. Note that altogether with  $MV$  and  $MTBV$ , these characteristics form the set from which proxy TS factors of the 5 Factor Model (Fama and French, 2015) are based on. The same winsorization scheme is applied for

Range	Full	US	Europe	BICS
1990-1994	2.62%	11.06%	1.70%	52.53%
1994-1999	1.01%	2.11%	1.16%	22.51%
2000-2004	10.39%	24.18%	7.54%	25.53%
2005-2009	5.71%	10.58%	2.53%	32.28%
2010-2014	3.93%	10.27%	3.58%	7.71%
2015-2019	1.97%	2.85%	3.03%	9.60%
2019-2024	1.09%	1.36%	3.66%	15.84%
1990-2024	29.59%	79.01%	25.49%	324.34%

TABLE 4.3: Cumulative returns by period and for the full sample - the range covers the period January 1990 to May 2024. Included variables: MTBV

Range	Full	US	Europe	BICS
1990-1994	2.54%	9.24%	1.75%	9.99%
1994-1999	-0.89%	-2.71%	0.24%	11.09%
2000-2004	10.48%	20.76%	8.70%	28.84%
2005-2009	5.45%	9.22%	3.00%	36.49%
2010-2014	3.55%	6.26%	3.38%	9.50%
2015-2019	2.01%	2.28%	4.35%	10.54%
2019-2024	1.34%	1.38%	3.34%	18.72%
1990-2024	26.73%	54.45%	27.29%	208.74%

TABLE 4.4: Cumulative returns by period and for the full sample - the range covers the period January 1990 to May 2024. Variables included: MTBV, OP, IG

each characteristic. In this setting  $F_{it} = \{MV_{it}, MTBV_{it}, IG_{it}, OP_{it}\}$  and  $X_{it} = \{\emptyset\}$ . Figure 4.8 plots the coverage of the dataset, considering those companies which have P, MV, MTBV, IG and OP available simultaneously. It is evident that we had to drop some observations at the end of the sample (the last available observations are on Dec 23) due to lack of data availability.

Table 4.4 shows very little impact in the size factor from the inclusion of IG and OP. The only exception to this statement are the returns for the Emerging Market, declining from 300% return over the full period to nearly 200%.

Finally, we included other control variables of interest: Price to Earnings (PE), Leverage (LEV) and Dividend Yield (DY). In this scenario  $F_{it} = \{MV_{it}, MTBV_{it}, IG_{it}, OP_{it}\}$  and  $X_{it} = \{PE_{it}, LEV_{it}, DY_{it}\}$ . 4.5 shows the returns for this specification, when adding all controls, the case of BICS is not available due to the limited number of data, so we focus on the full dataset, EU and US case.

Range	Full	US	Europe
1990-1994	1.01%	4.05%	0.68%
1994-1999	-3.06%	-7.46%	-0.56%
2000-2004	8.12%	10.97%	9.90%
2005-2009	3.98%	4.61%	2.83%
2010-2014	3.61%	5.28%	4.17%
2015-2019	1.10%	-0.19%	3.92%
2019-2024	1.15%	-0.53%	2.45%
1990-2024	16.63%	16.82%	25.48%

TABLE 4.5: Cumulative returns by period and for the full sample - the range covers the period January 1990 to May 2024. Variables included: MTBV, OP, INV, PE, LEV, DY

The variations in results obtained from different model specifications can be attributed to the (cross-sectional) correlations between the observed characteristics considered.

Figure 4.15 provides a visual representation of the correlation between Market Value (MV) and the other characteristics included in the model at each point in time. This visualization allows for a more nuanced understanding of how these characteristics interact over the observed period and highlights any temporal trends that may affect the returns resulting from the cross-sectional regressions.

As mentioned previously, the magnitude of these correlations remains relatively low throughout the analyzed period, with the Market-to-Book Value (MTBV) ratio emerging as the variable most strongly correlated with MV, particularly during the years spanning from 1990 to 2005. This suggests that the differences between cumulative returns in the reported figures may be somewhat “amplified” by the compounding nature of cumulative returns, as even small differences in the returns could have a substantial impact on the cumulative returns illustrated.

This pattern raises several questions regarding the purity and interpretability of the estimated effects, as portfolio returns obtained from this procedure may be affected to what is known as “bias from omitted variables.” Specifically, the omission of certain variables that are correlated with MV can result in biased estimates of the size effect. While the inclusion of additional characteristics and control variables provides a partial solution to this issue, it also leads to a reduction in dataset coverage, thus decreasing the sample size and potentially impacting the robustness of the results.

Moreover, it is plausible that latent, unobserved factors exert an influence on the results, given the complexity of cross-sectional data. This possibility, however, remains

largely unexplored in the existing literature on cross-sectional factors, which often focuses on observable characteristics without addressing potential underlying dimensions that may further contextualize or modify observed effects. Future research could benefit from incorporating methods to account for latent variables, thus offering a more comprehensive and robust methodology to retrieve cross-sectional factors.

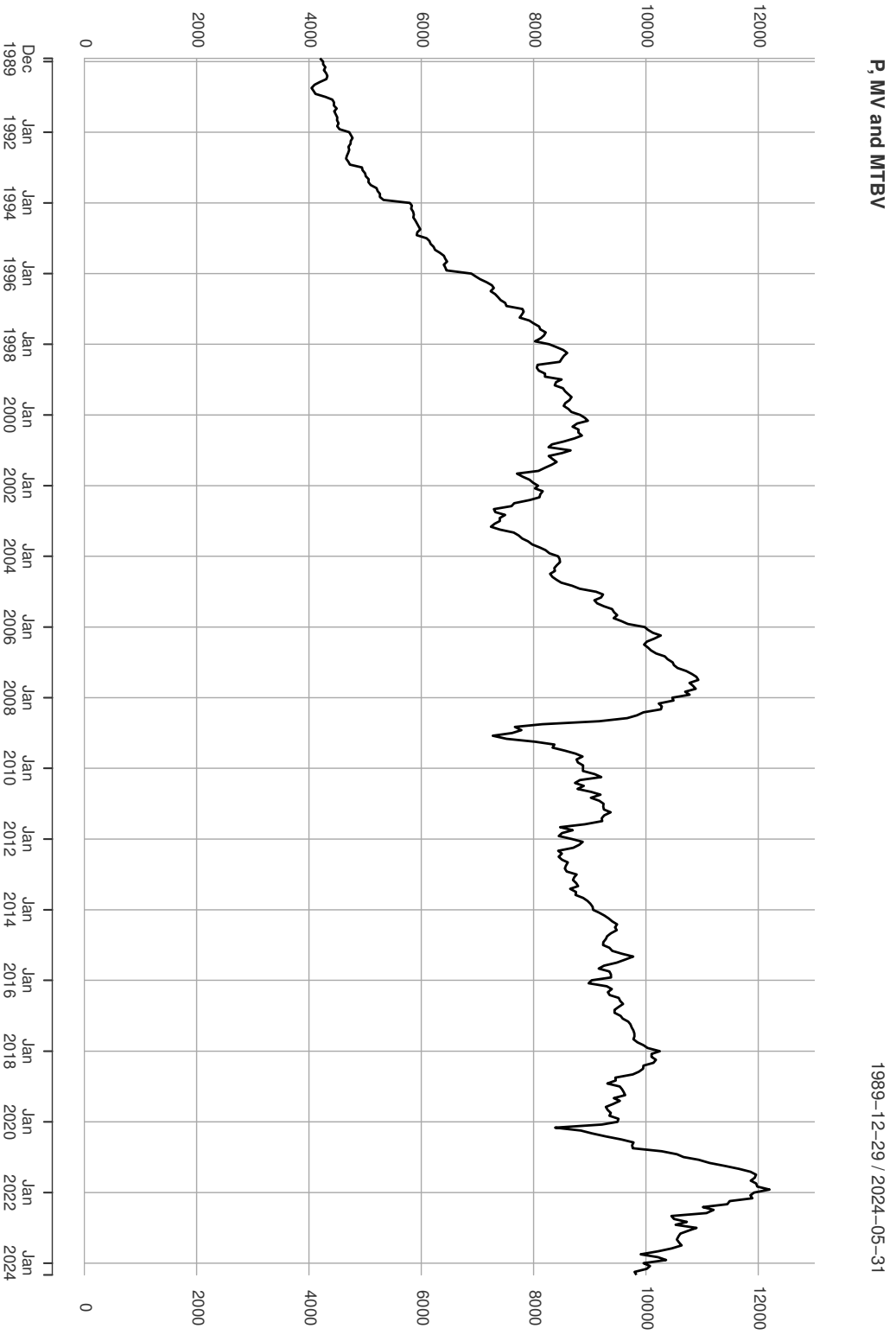


FIGURE 4.7: Dataset coverage, number of companies for each month, when Price (P), Market Value (MV), and Market-to-Book value (MTBV) are all available.

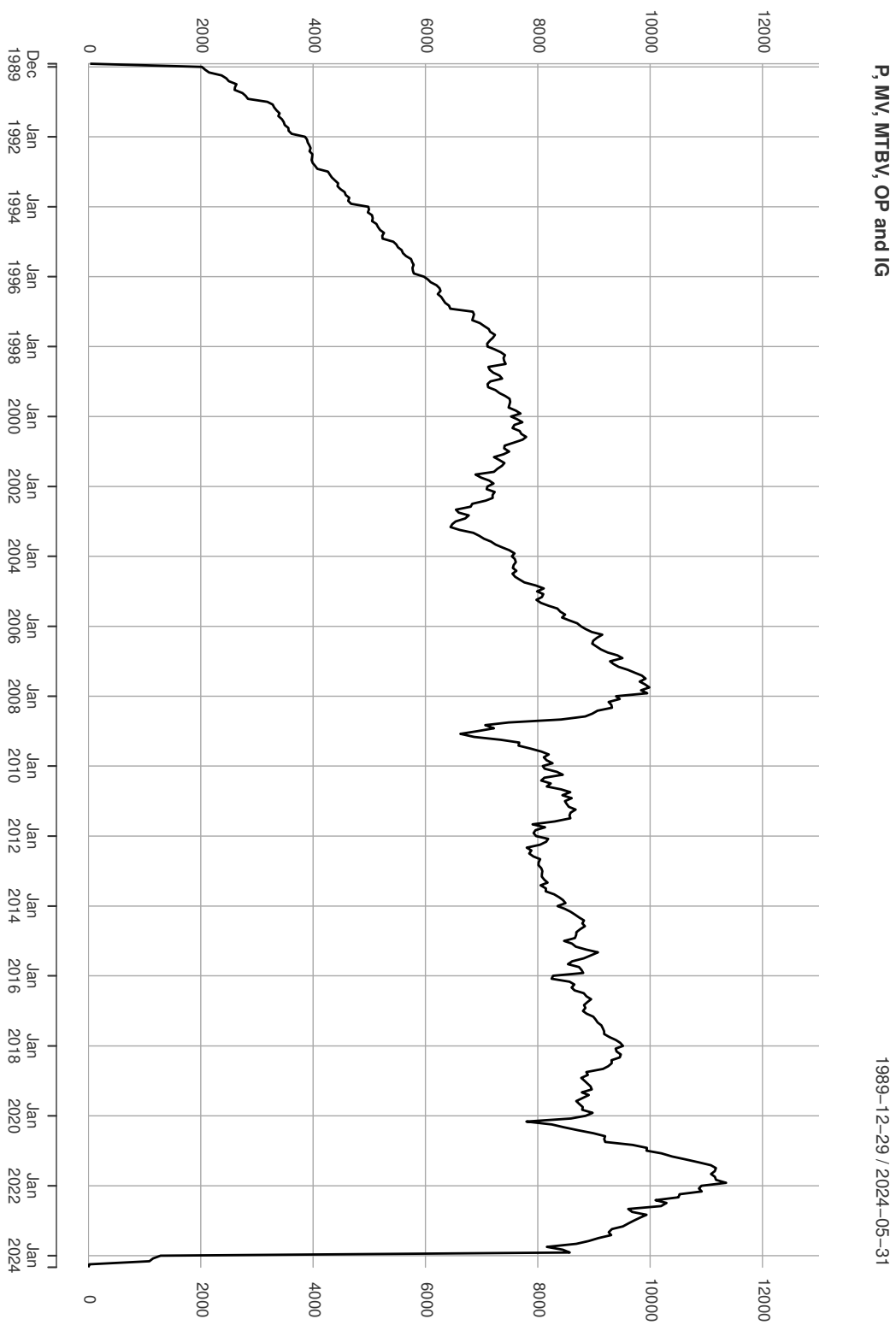


FIGURE 4.8: Dataset coverage, number of companies for each month, when Price (P), Market Value (MV), Market-to-Book value (MTBV), Operating Profit (OP), and Investment Growth (IG) are all available.

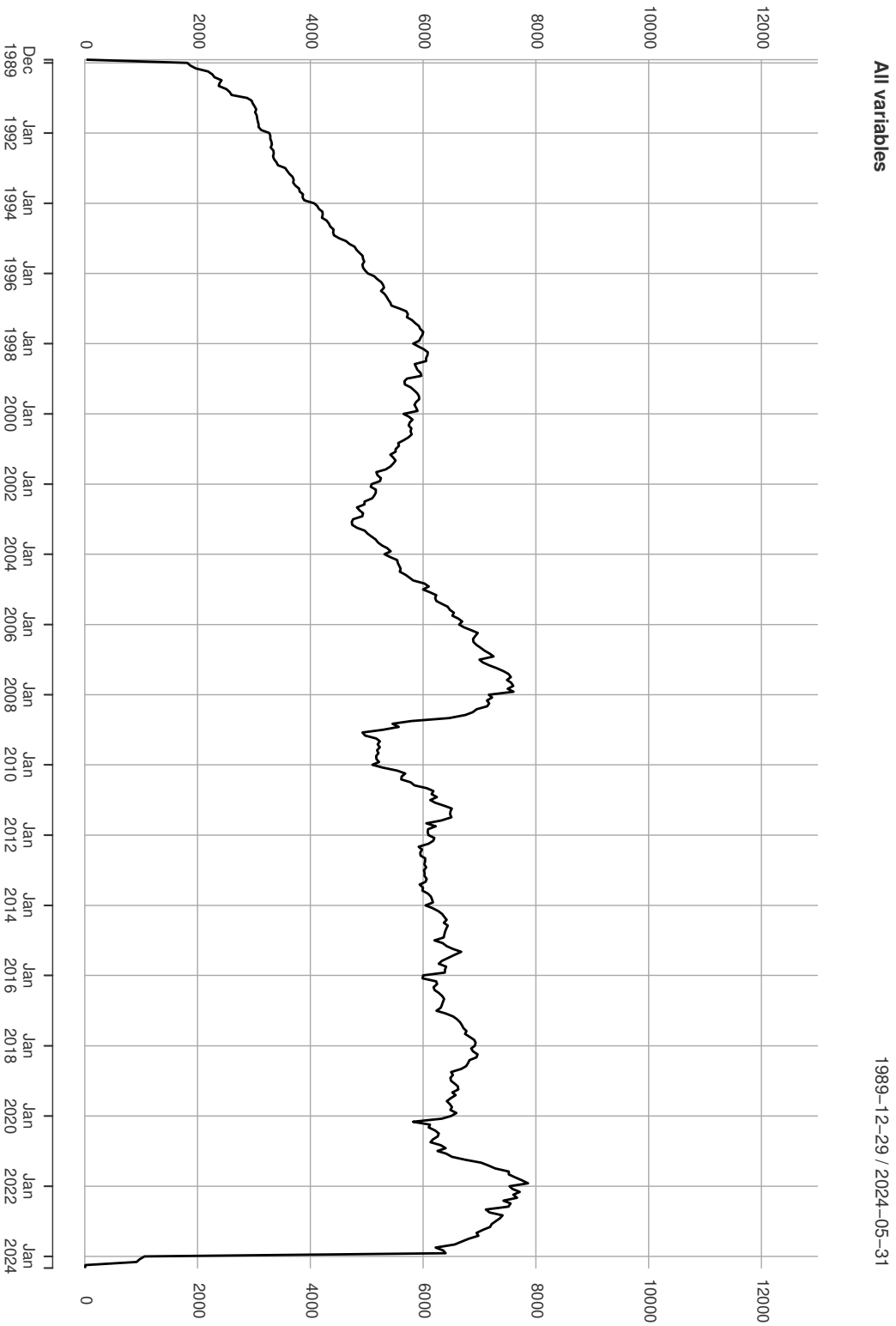


FIGURE 4.9: Dataset coverage, number of companies for each month, when Price (P), Market Value (MV), Market-to-Book value (MTBV), Operating Profit (OP), and Investment Growth (IG), Price/Earnings (PE), Leverage (LEV), and Dividend Yield (DY), are all available.

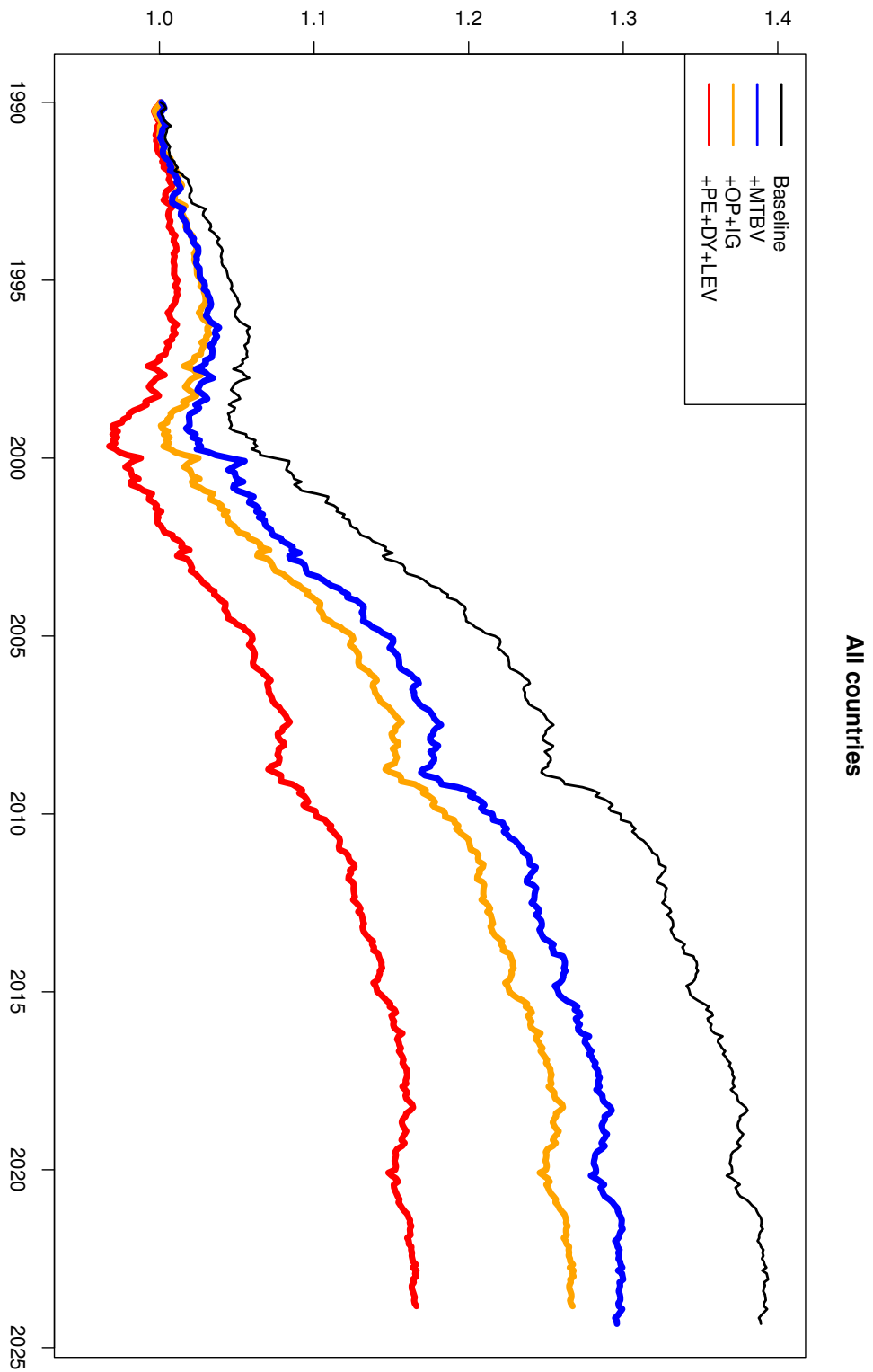


FIGURE 4.10: Comparison of alternative designs for the Size factor estimation. Cumulative returns. Data from 1990 - All countries

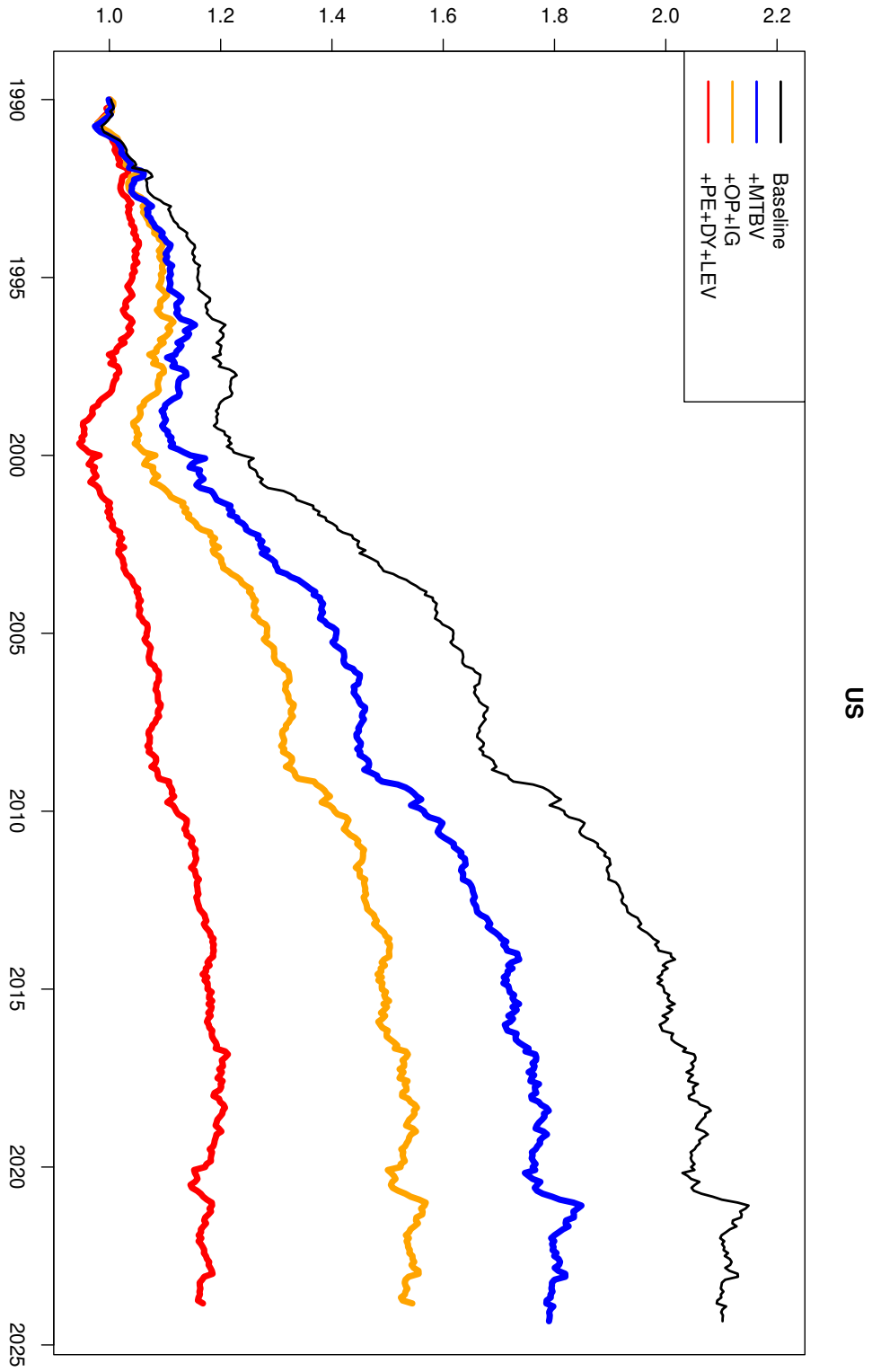


FIGURE 4.11: Comparison of alternative designs for the Size factor estimation. Cumulative returns. Data from 1990 - United States markets.

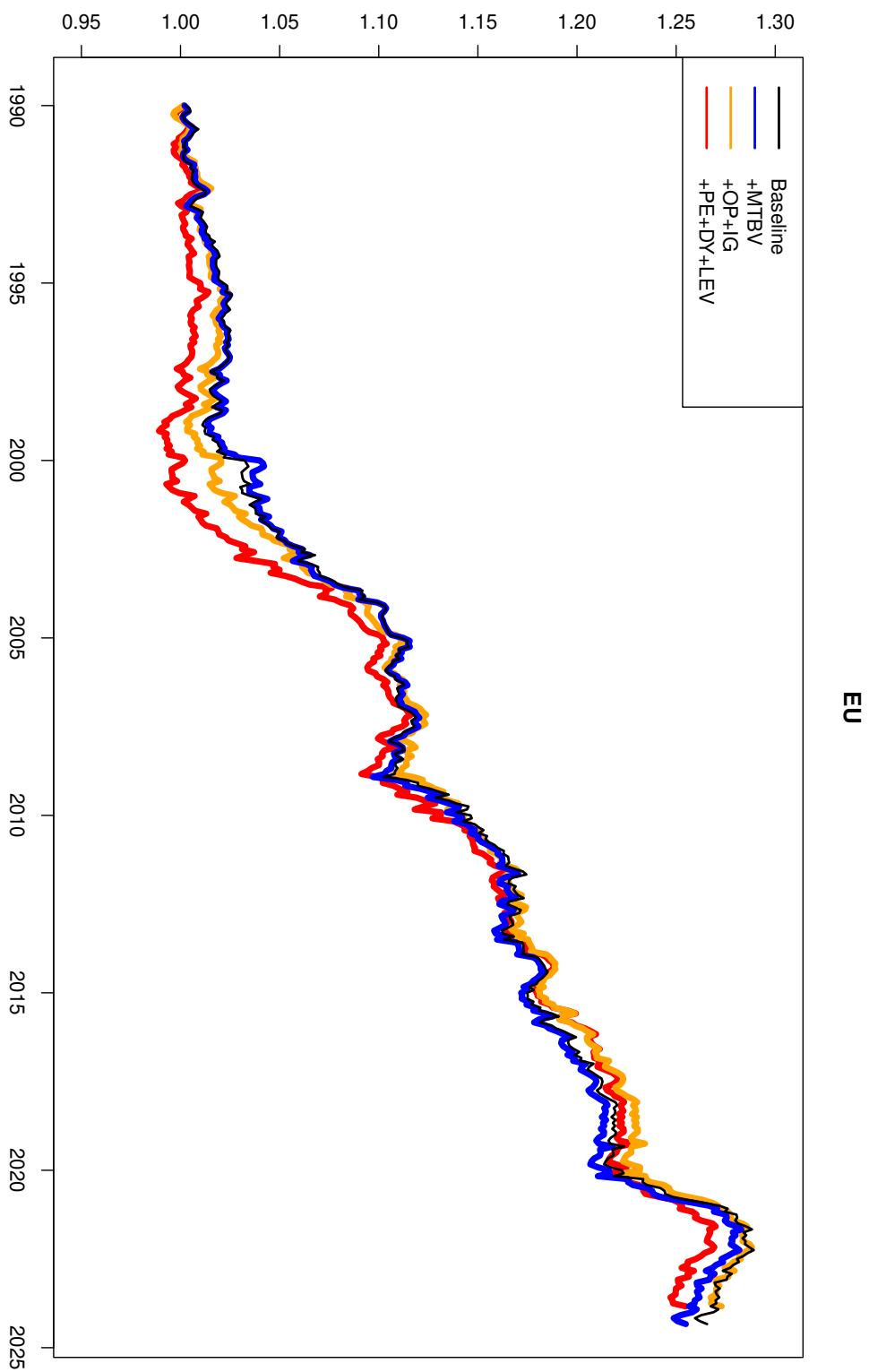


FIGURE 4.12: Comparison of alternative designs for the Size factor estimation. Cumulative returns. Data from 1990 - European markets

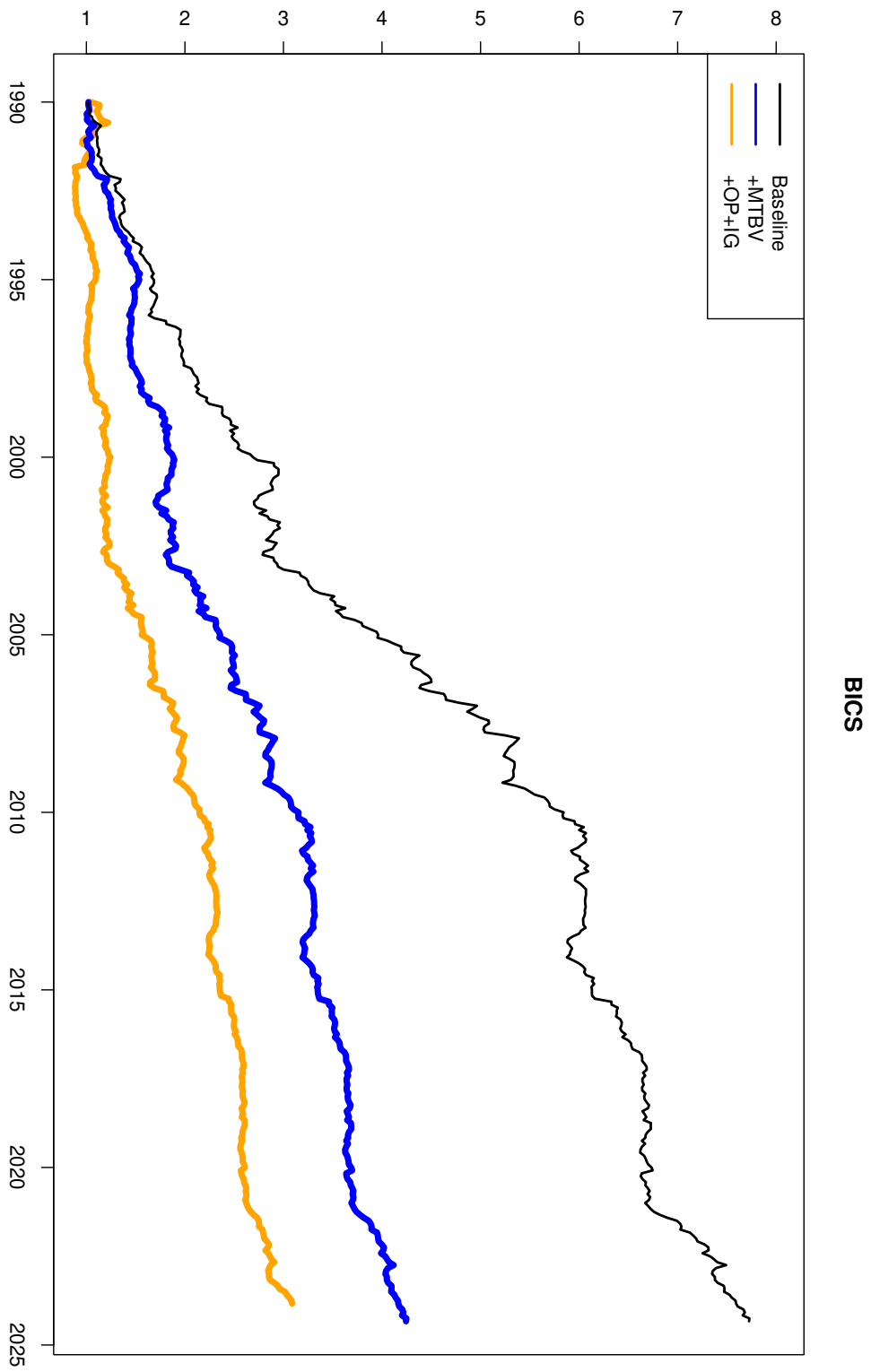


FIGURE 4.13: Comparison of alternative designs for the Size factor estimation. Cumulative returns. Data from 1990 - Emerging markets

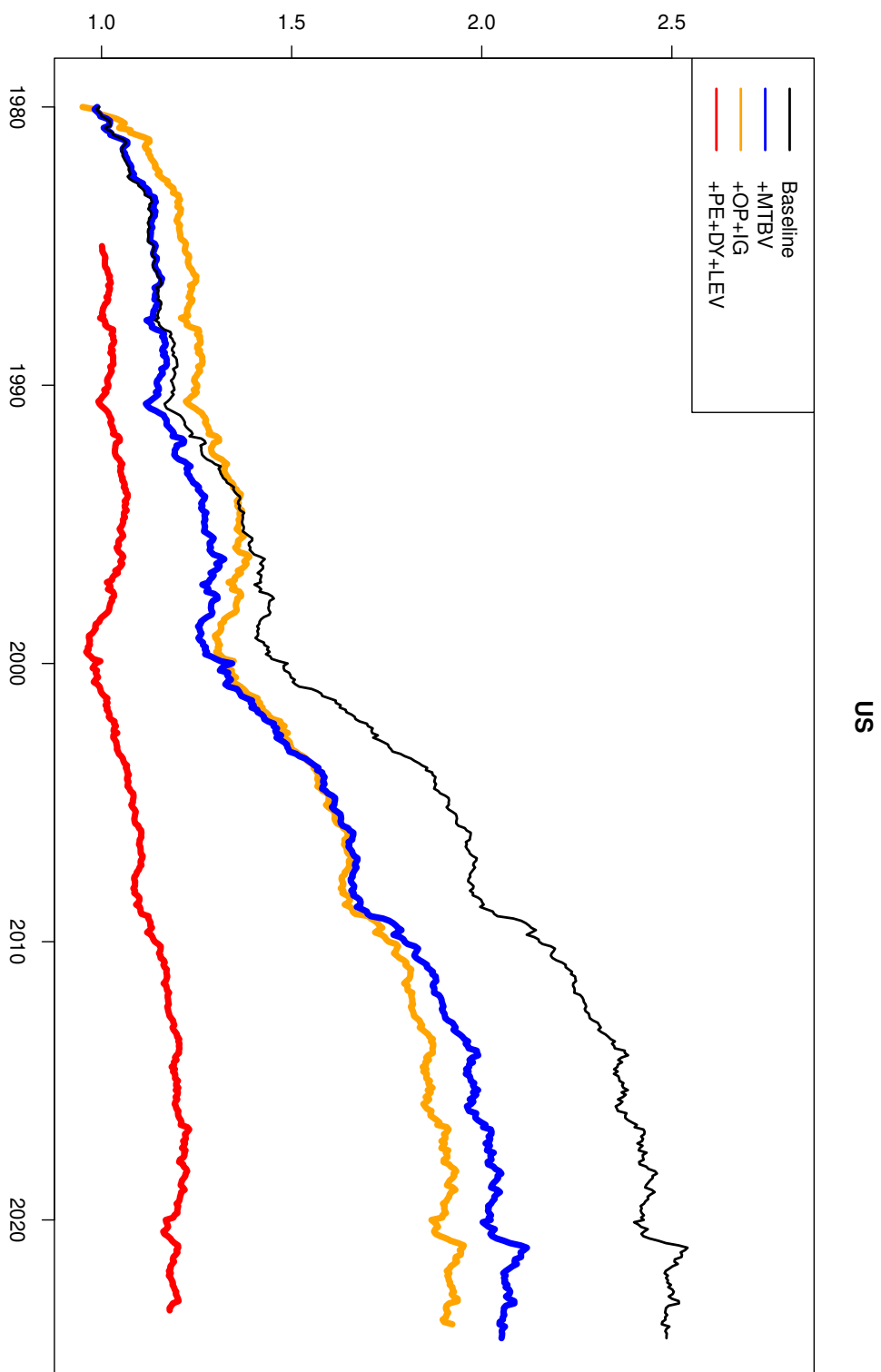


FIGURE 4.14: Comparison of alternative designs for the Size factor estimation. Cumulative returns. Data from 1980 - United States markets.

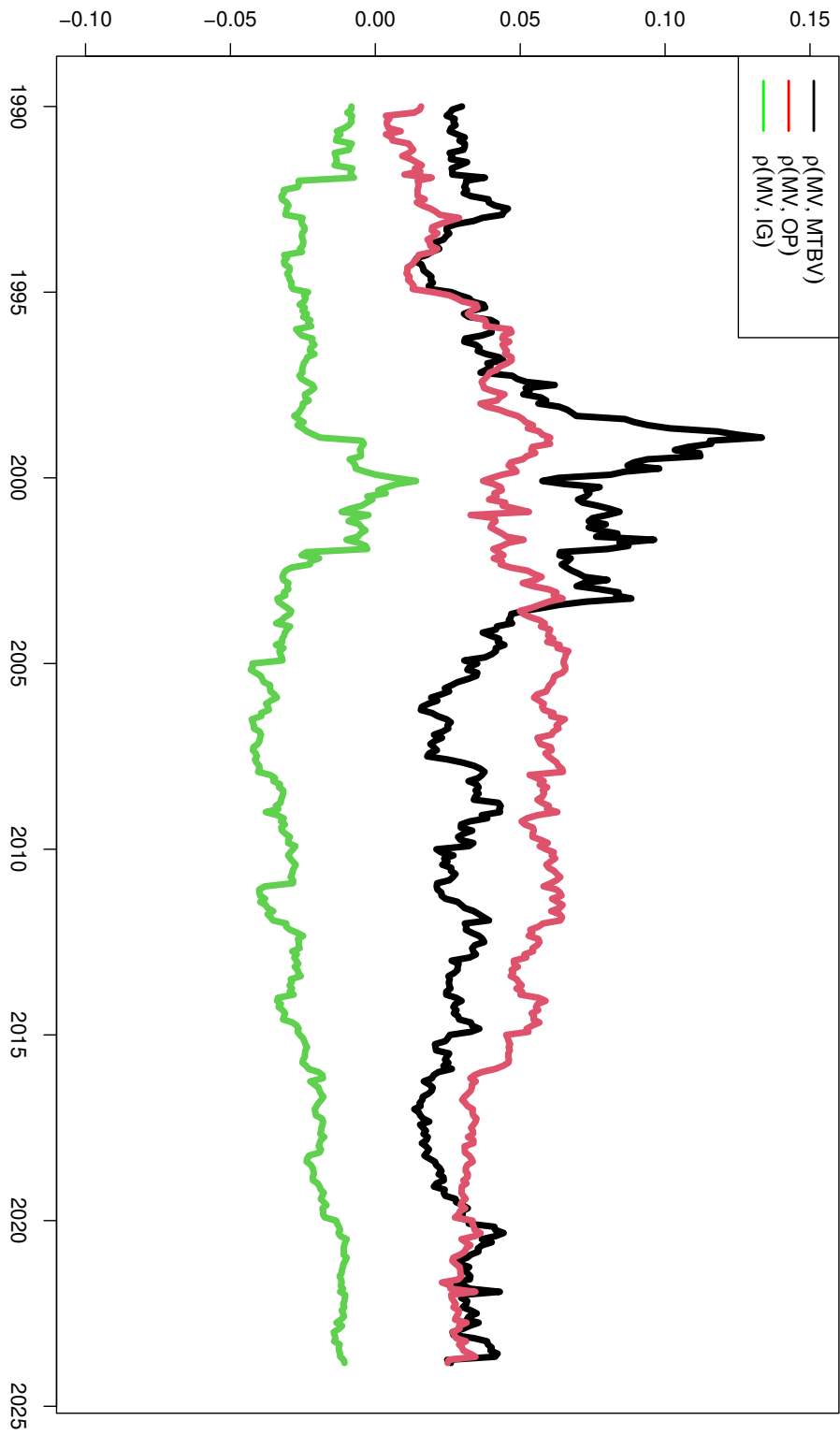


FIGURE 4.15: Cross-sectional correlation between MV and other characteristics at each point in time. Data from 1990 - Full Sample

### 4.3.2 The role of Megacaps

In this section, we investigate the influence of the largest companies (in terms of Market Value) within the previously derived size factor. Given the systematic exclusion of penny stocks from our dataset, we anticipate a skewed distribution relative to company size, likely characterized by a pronounced right tail. This distributional skewness introduces certain limitations, especially given our chosen methodology. Specifically, the portfolio weights in the cross-sectional size factor are derived from standardized characteristics across observations. As a result, variations in the distribution of securities, particularly by Market Value, inherently alter the portfolio's weight composition.

A further question that arises concerns whether the recent outperformance of large-cap firms indicates a potential shift away from the small-cap effect, suggesting instead a “large-cap effect” that may persist for a number of years. This question holds particular relevance for the so-called “Magnificent 7” companies, namely: Apple, Microsoft, Amazon, Alphabet, Tesla, Meta, and Nvidia.

Figures 4.17 and 4.18 illustrate the cumulative size factor obtained from excluding the Magnificent 7 companies, alongside the baseline scenario for the full dataset and its US sub-sample. While the results are not directly comparable due to differing datasets, which influence the standardized values and, consequently, the weights of commonly included companies, certain notable patterns can still be observed. Naturally, the impact of excluding these companies is more pronounced in the US sub-sample respectively; however, this effect may partially stem from the relatively higher proportion of observations excluded within the US sample compared to the global dataset.

In both cases, we find that excluding the Magnificent 7 results in higher cumulative returns. This divergence appears around 2015 for the full sample, while it becomes evident earlier, around 2008, in the US sub-sample. The absence of any observable effect in the early years of the dataset could reflect the recent performance and significant growth of these companies in recent years. Nevertheless, all of these effects, as we observe, manifest merely as contractions (or expansions when excluding the companies) in cumulative returns, without any indication of a change in the underlying return pattern.

The impact of the “Magnificent 7” can be extended by broadening the focus to include a wider group of large-cap companies, allowing for a dynamic selection of excluded companies that varies over time rather than remaining fixed. Figure 4.19 presents the results of excluding top companies (ranked by Market Value) at varying thresholds.

The primary observation is that, as more companies are excluded, cumulative returns

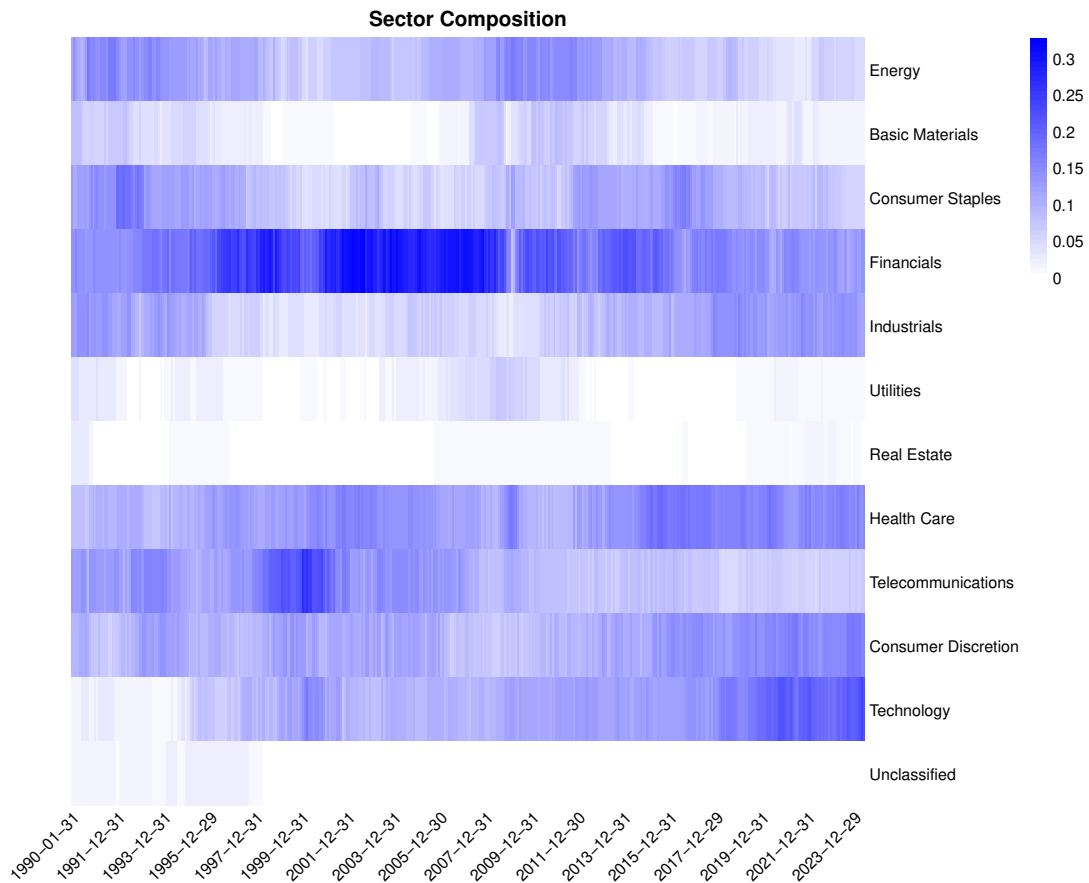


FIGURE 4.16: Sector composition of the last percentile of assets by MV in each cross-section

increase. Notably, the onset of this effect occurs earlier than in the case of the Magnificent 7 exclusion. This finding suggests that the observed reduction in the effectiveness of the small-cap effect cannot be solely attributed to the recent outperformance of the Magnificent 7. If such outperformance does play a role, it likely accounts for only a portion of this shift.

Again, it is important to emphasize that each exclusion threshold results in a distinct dataset, thus yielding different standardized characteristics and portfolio weights across securities. This variability prevents a precise “quantitative” comparison across the thresholds used, yet it offers valuable insight into how the exclusion of the largest companies might impact a potential small-cap strategy. Overall, while the exclusion of large caps influences cumulative returns, it does not alter the underlying pattern of returns, highlighting a consistent behavioral trend across these scenarios.

Figure 4.16 illustrates the changing sector composition within the top 1% of assets ranked by market value (MV) over time. This elite subset represents large-cap assets that typically exert a downward effect on the returns of the size factor. A closer examination of the sectoral distribution within this group thus provides insights into structural shifts affecting large-cap dynamics.

The *Financials* sector is the most consistently represented, accounting for nearly 30% of the assets under study from 1995 through 2007. It is also noteworthy that the *Technology* sector has gained increasing prominence in recent years, a trend exemplified by the rise of the “Magnificent 7” within this sector. In contrast, sectors such as *Basic Materials* and *Utilities* are entirely absent from this high-value subset.

Figures 4.20 and 4.21 present the median values of various company-specific characteristics over time for companies grouped by different percentiles based on Market Value (MV). Specifically, we examine the first decile, the last decile, and the top 1%. This comparison provides insights into how small-cap companies differ from large-cap companies.

Our primary focus here is on the differences between the top 1% of companies and those in the last decile, aiming to determine whether any fundamental characteristics drive the out-performance of the largest-cap companies. Overall, we observe that most characteristics show no significant differences between these groups. However, the top 1% of companies exhibit a higher Market-to-Book Value (MTBV), with this gap widening over the past decade. A similar pattern is observed for Earnings per Share (EPS), and these companies also display higher Dividend Yields.

These differences, however, appear insufficient to fully explain the return disparities in the size factor observed when the largest companies are excluded. This finding suggests that factors beyond individual company characteristics may have contributed to the recent out-performance of large-cap firms. For instance, recent discussions in financial literature have highlighted concerns that passive investing may be fueling a bubble in large-cap stocks. This phenomenon is primarily attributed to the structure of passive investment vehicles, such as exchange-traded funds (ETFs) and index funds, which are typically market-capitalization-weighted. In these funds, larger companies constitute a more significant portion of the portfolio, leading to a disproportionate allocation of capital to large-cap stocks. Critics argue that as more investors channel funds into passive strategies, the demand for large-cap stocks increases, potentially inflating their valuations beyond fundamental justifications.

This is an aspect worth considering, especially in light of our analysis results; however, it falls beyond the scope of this work.

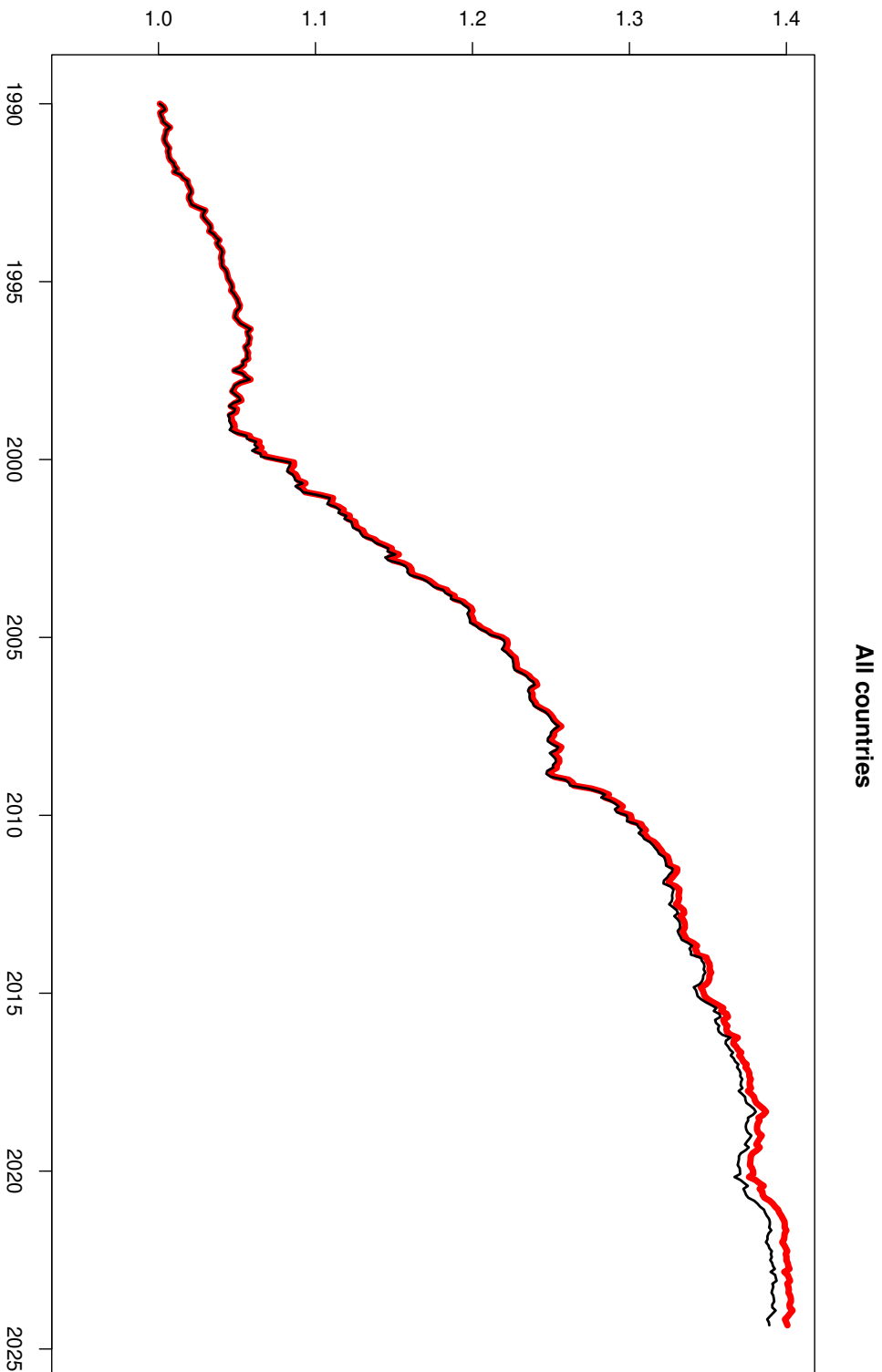


FIGURE 4.17: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using only the Market Value as explanatory variable and excluding the *Magnificent 7* - data from 1990 - entire dataset. Black line: estimate without excluding the *Magnificent 7*.

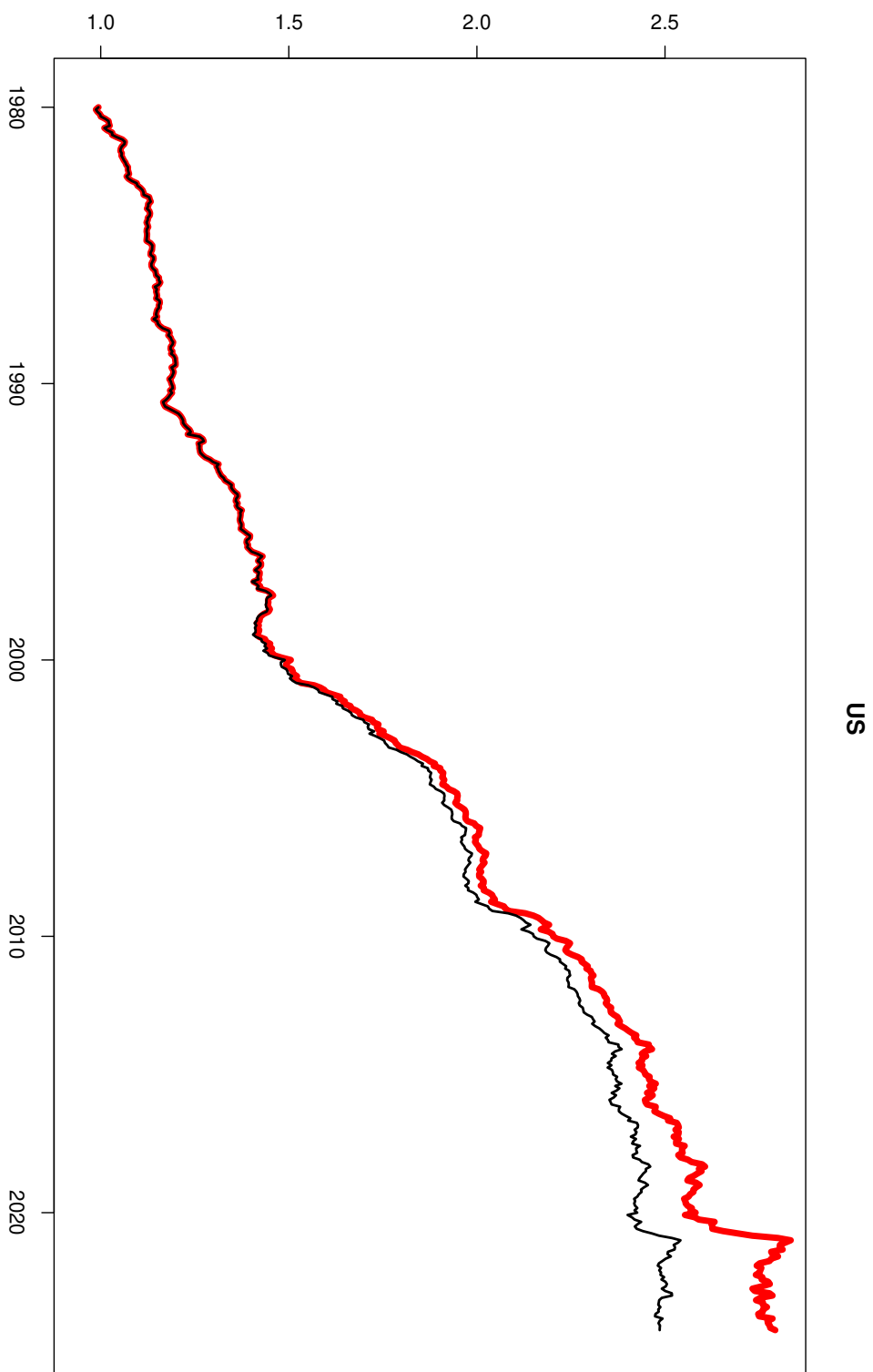


FIGURE 4.18: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using only the Market Value as explanatory variable and excluding the *Magnificent 7* - data from 1980 - US markets. Black line: estimate without excluding the *Magnificent 7*.

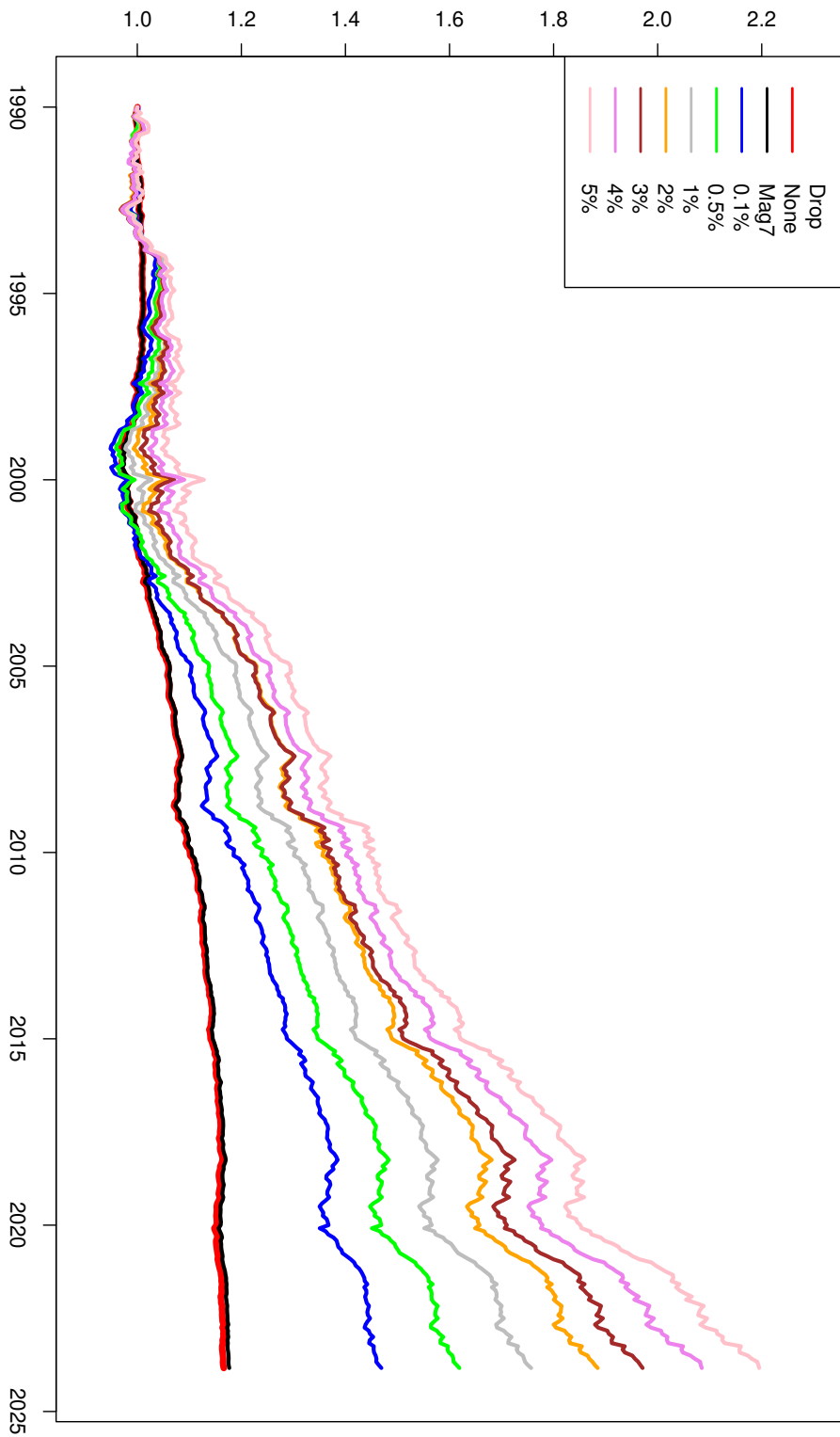


FIGURE 4.19: Cumulative returns of the Size factor based on cross-sectional estimates using only the full specification and excluding the top companies by Market Value with different relative thresholds - data from 1990 - all countries.

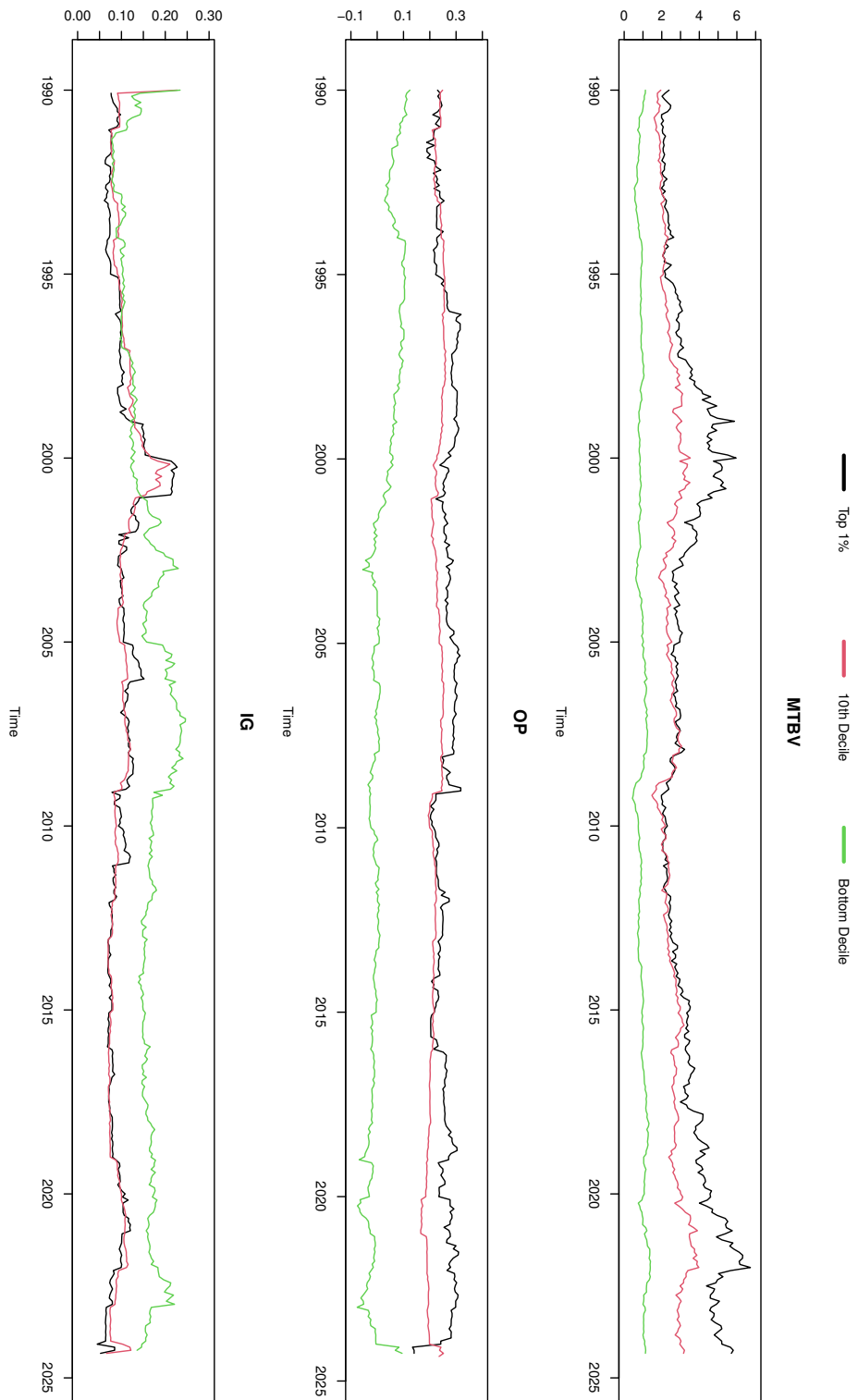


FIGURE 4.20: Median Values of Characteristics for Companies Across Different MV Percentiles

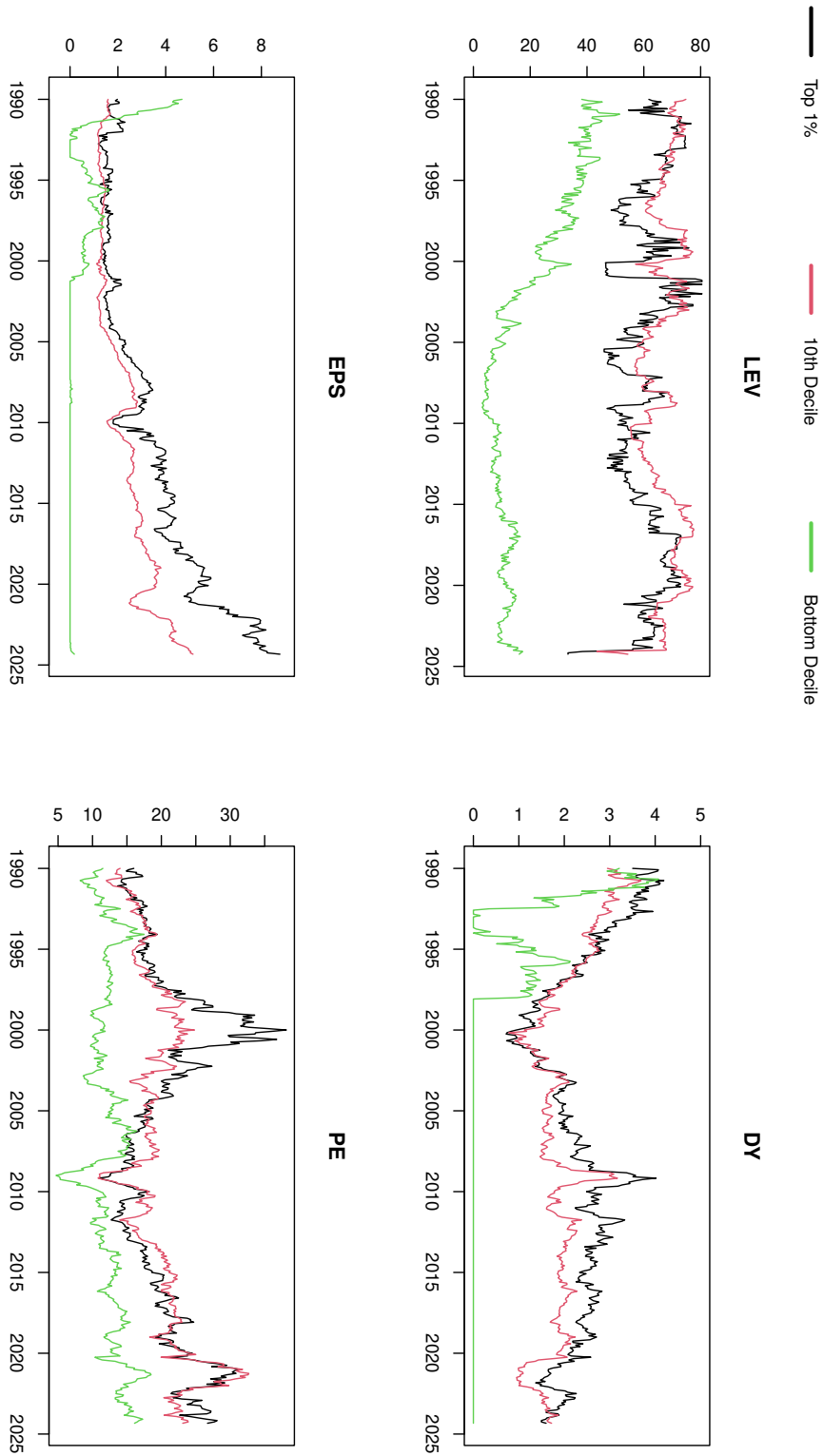


FIGURE 4.21: Median Values of Characteristics for Companies Across Different MV Percentiles

### 4.3.3 Other drivers of the size factor

Building on the findings from the previous section, where the exclusion of the largest companies resulted in a pronounced shift in the size factor, we are led to consider whether other financial characteristics might also significantly influence the returns associated with this factor. This question guides us to explore additional drivers beyond firm size that could be shaping the dynamics of the size factor over time.

To investigate this, we apply a similar approach, recalculating the size factor by systematically excluding observations that fall within the highest percentiles of the distribution for specific financial characteristics. This strategy allows us to derive alternative estimates of the size factor that reflect the potential influence of each characteristic independently.

The financial variables analyzed in this context include the Market-to-Book Value (MTBV), Leverage (LEV), Earnings per Share (EPS), and the inverse of the Price-to-Earnings (PE) ratio—effectively capturing the Earnings-to-Price (1/PE) measure.

The results, as shown in Figure 4.22, illustrate the differentials between the standard size factor, computed with MTBV included as an explanatory variable in the cross-sectional regressions (and incorporating the usual winsorization scheme), and the modified size factor obtained after removing companies within the top 5% for each characteristic, still including MTBV. It is important to note that for the MTBV characteristic itself, the exclusion of the top 5% subset occurs post-winsorization, ensuring consistency in the treatment of outliers. The differentials observed serve as proxies for the extent to which each variable contributes to the size factor.

Our findings reveal that these characteristics do, in fact, impact the size factor, with Earnings-to-Price (1/PE) exhibiting the most substantial effect among the variables examined. Notably, the differentials remain (almost) consistently positive throughout the entire period under study, indicating that, unlike the case with large-cap exclusions, the removal of assets with high values in these characteristics tends to result in a contraction of the cumulative size factor.

However, it is critical to observe that the influence exerted by these financial characteristics, while significant, is remarkably smaller in magnitude compared to the effect of excluding large-cap companies. This contrast highlights the dominant role that large-cap exclusions play in evaluating the size factor, reinforcing the intuition that excluding large-cap firms may be crucial to implement a robust small-cap strategy.

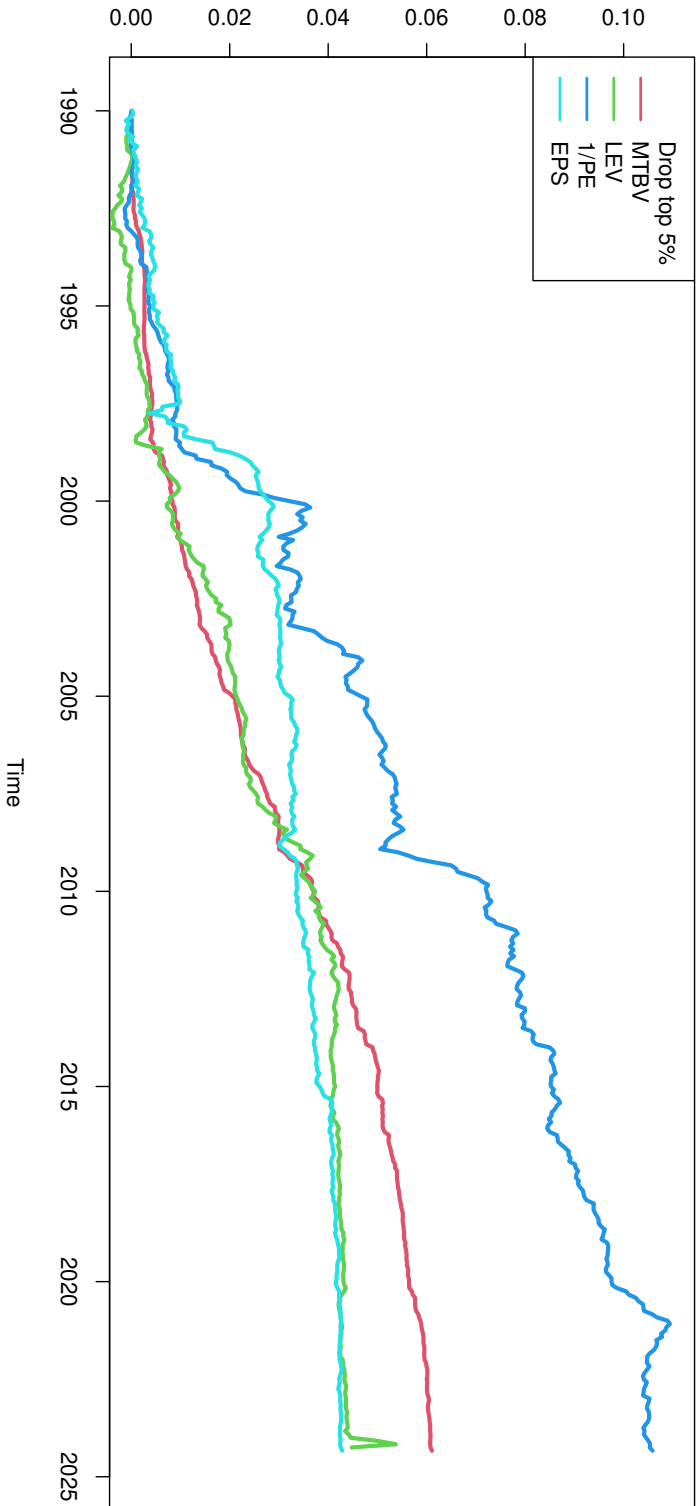


FIGURE 4.22: Differentials between size factor for the full dataset when MTBV is included as explanatory variable in the cross-section regression and size factor when the top 5% of company with respect to other variables are excluded

# Chapter 5

## Time-Series analysis

In this chapter, we transition to a time-series framework, applying a time-series regression for each individual firm. The baseline model utilized is the Fama-French three-factor model, where the dependent variable is the excess market return of each company. The risk-free rate used in the analysis is the 1-month Treasury bill. The explanatory variables comprise the market return, proxied by the MSCI All Country World Index returns, and two factors related to firm size and the market-to-book ratio (SMB and HML).

Notably, these factors are replaced by the cross-sectionally estimated factors derived in the preceding chapter, where the market-to-book value (MTBV) ratio is winsorized to mitigate the influence of extreme values. Additionally, we will evaluate the strength of the factors, following the methodology outlined in Chapter 1.

### 5.1 Time-Series results

This section presents the primary results from estimating a time-series asset pricing model, specifically focusing on the three-factor model originally introduced by Fama and French. As the model name suggests, we employ three factors: the market excess return, the Size factor, and the Value factor. Here, the SMB (Size) and HML (Value) factors proposed by Fama and French are replaced by their cross-sectional counterparts, derived in the preceding chapter. Specifically, we utilize the factors obtained by including only Market Value (MV) and Market-to-Book Value (MTBV) as explanatory variables in the cross-sectional regressions.

Ideally, the objective would be to use a size factor that controls for the largest possible number of characteristics, resulting in a “purer” factor, or one that minimizes bias due to omitted variables. However, we have chosen to restrict the model to MV and MTBV,

primarily for reasons of data availability. As detailed in the previous chapter, introducing additional variables reduces data coverage due to availability constraints. Moreover, MTBV demonstrated a considerable impact on the resulting size factor, making this specification a pragmatic “trade-off” between factor purity and data coverage.

The model specification is as follows:

$$R_{it} - R_{ft} = \alpha_i + \beta_{1i}(R_{mt} - R_{ft}) + \beta_{2i}R_{MCt} + \beta_{3i}R_{MTBVt} \quad (5.1)$$

This model is estimated for each individual firm, producing an estimated regression coefficient for each factor for every firm in the dataset.

In this section, we provide the main findings for two scenarios: one using the full time period and the other employing 10-year rolling windows, with a 6-month step between consecutive windows. It is worth noting that our primary interest does not lie in assessing model performance or evaluating its predictive power, as the focus remains on analyzing the size factor. Nonetheless, the model serves as a foundational element for assessing factor strength, which will be explored in the following section.

With the shift from cross-sectional to time-series regression, several methodological adjustments are necessary, differing from those previously employed. For consistency, we include only assets with a complete set of observations over the period considered, ensuring that differences in estimates are not attributable to varying sample sizes across periods. Additionally, penny stocks remain excluded, with the condition now being that a stock is classified as a penny stock if its price is 5 USd or below for at least half of the period considered.

Table 5.1 provides summary statistics for the intercepts across the entire sample period. These results indicate robust model performance, comparable to that reported in Chapter 2 of Fama and French’s original work.

$A \alpha $	$A t(\alpha) $	$As(\alpha)$	$As(\epsilon)$
0.007	1.859	0.006	0.115

TABLE 5.1: Summary statistics for estimated intercepts

To capture the temporal evolution of the estimated betas, we employ a rolling window approach. Figures 5.1 to 5.4 depict the median, first quartile, and third quartile of the estimated size factor betas, offering both geographical and sectoral perspectives. It is essential to note that this evaluation of betas is preliminary, focusing only on the estimated coefficients without accounting for their standard errors, and thus without

assessing precision. A more rigorous evaluation will be provided in the next section, where we analyze the factor strength to achieve a statistically robust understanding.

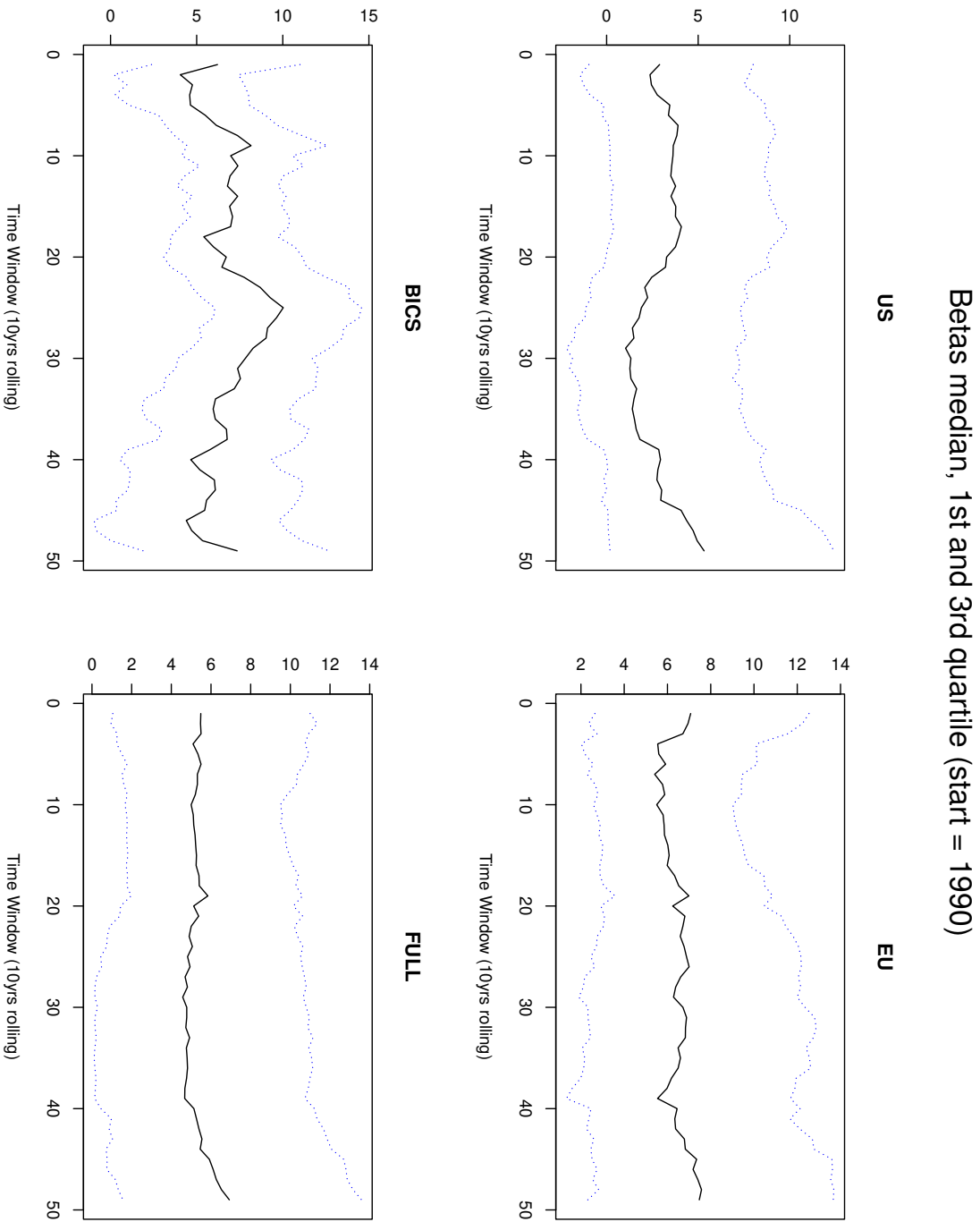


FIGURE 5.1 : Estimated exposures for the CS size factor by geographical area, 10 years rolling windows with 6m step (Full Sample from 1990).

Betas median, 1st and 3rd quartile (start = 1990)

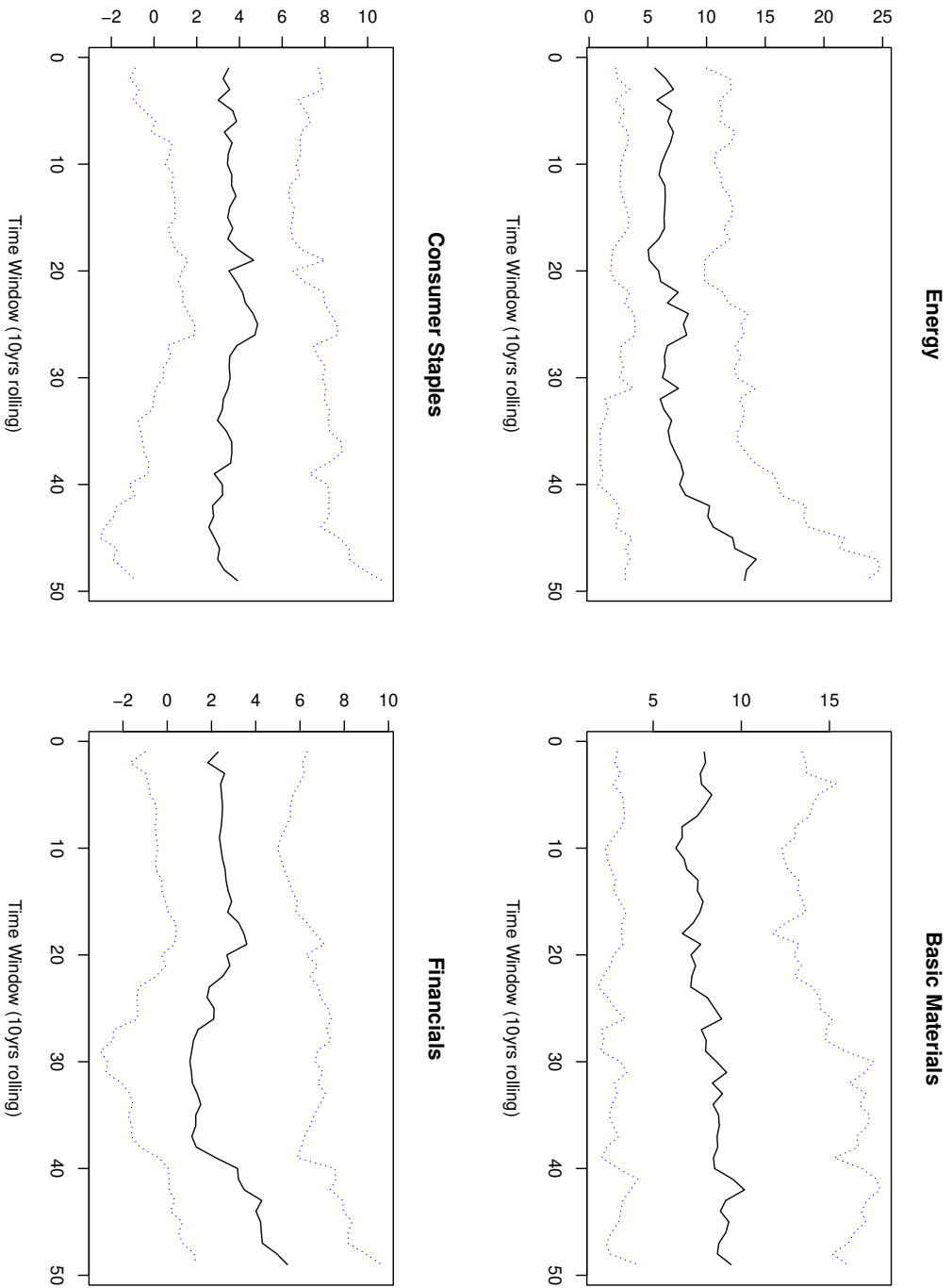


FIGURE 5.2: Estimated exposures for the CS size factor by sector, 10 years rolling windows with 6m step (Full Sample from 1990).

Betas median, 1st and 3rd quartile (start = 1990)

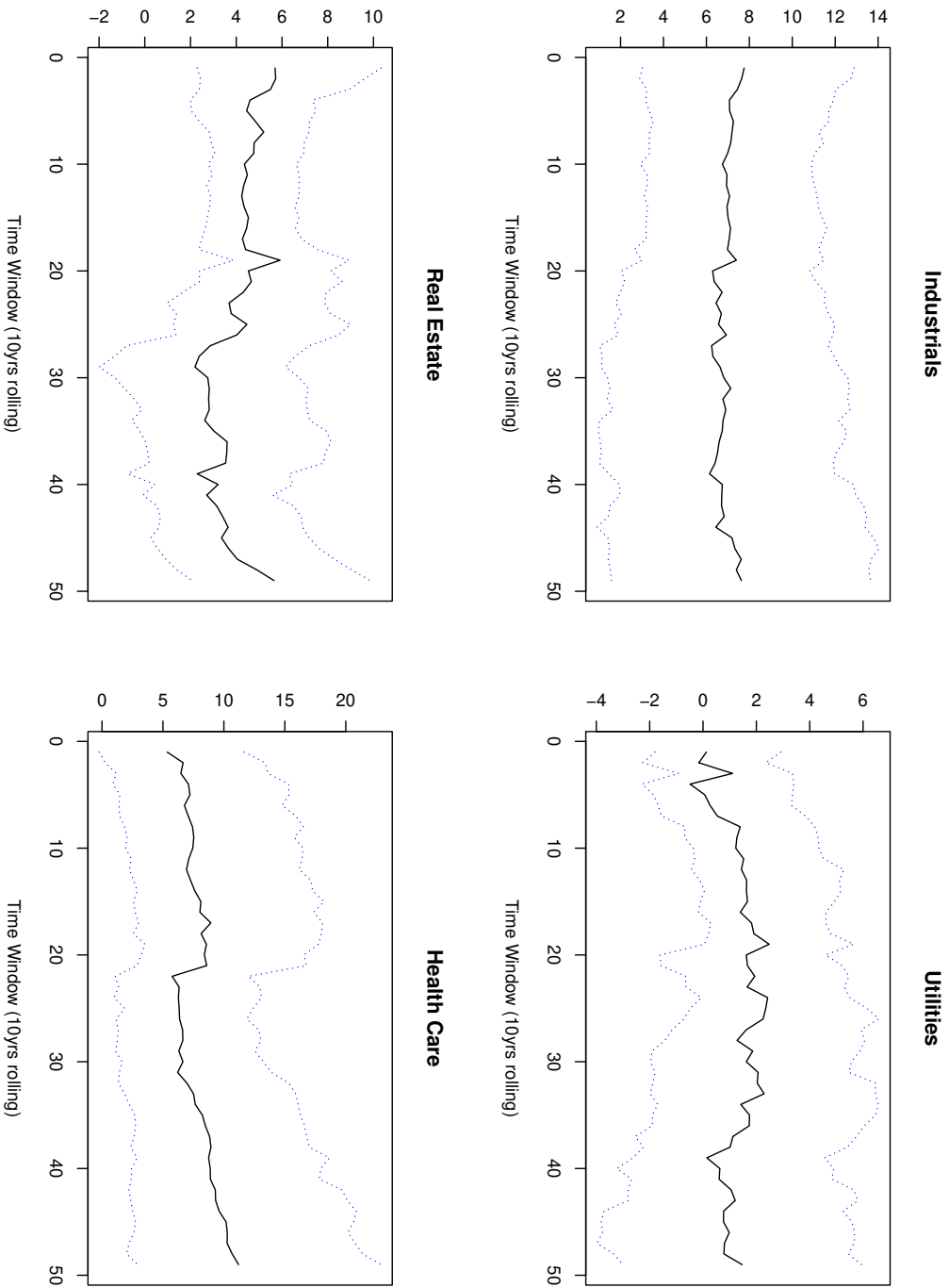


FIGURE 5.3: Estimated exposures for the CS size factor by sector, 10 years rolling windows with 6m step (Full Sample from 1990).

Betas median, 1st and 3rd quartile (start = 1990)

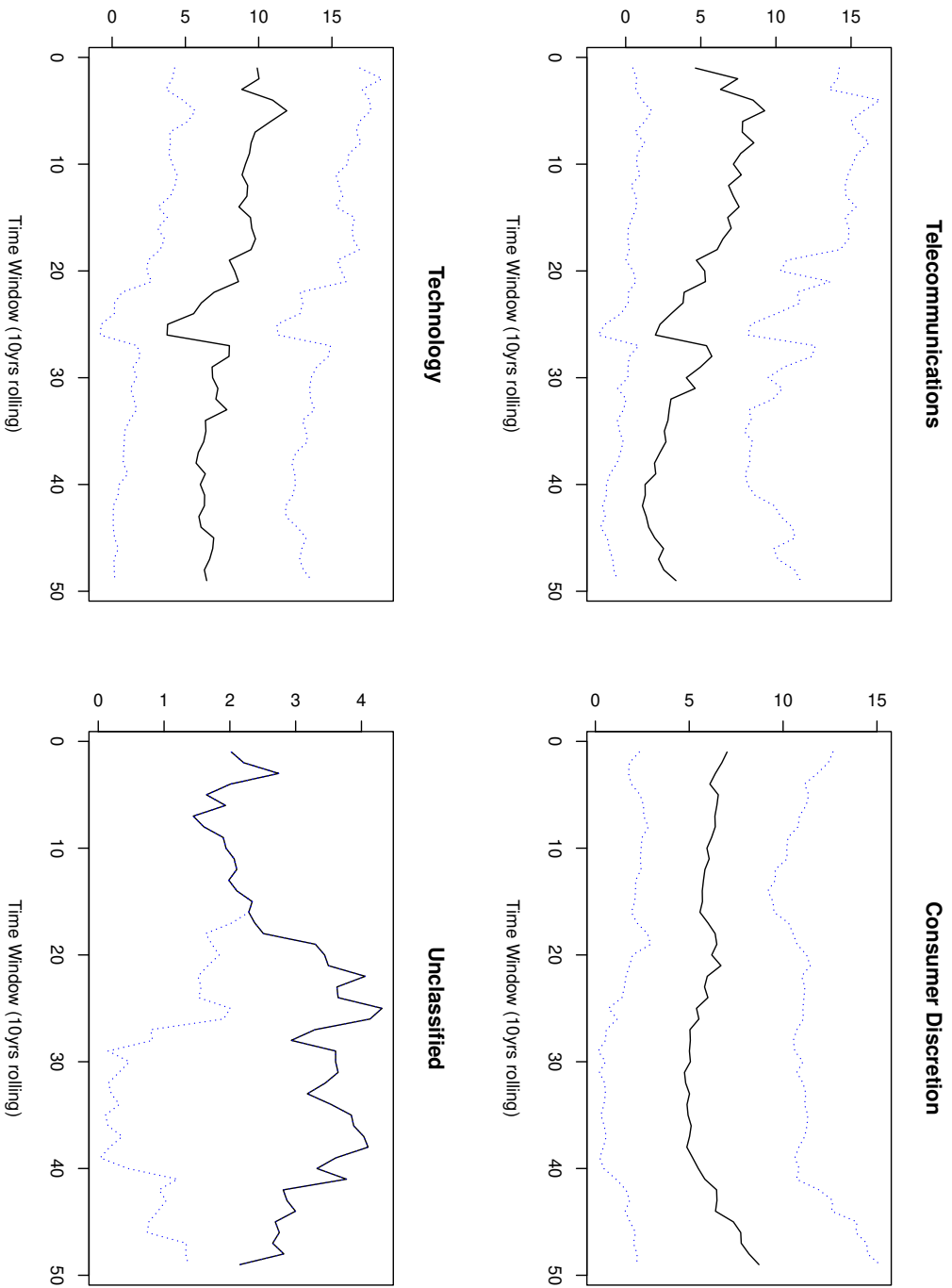


FIGURE 5.4: Estimated exposures for the CS size factor by sector, 10 years rolling windows with 6m step (Full Sample from 1990).

## 5.2 Evaluating Factor Strength

This section presents the results of the factor strength estimation. We have previously discussed how factor strength impacts the estimation and inference of risk premia using the two-pass procedure. However, in this context, we employ factor strength for a different purpose.

Since our methodology diverges from the two-pass approach, having already retrieved the monthly risk premia, our aim here is to assess the pervasiveness of the estimated factor in relation to the time series returns of securities, both at an aggregate and a local level.

Factor strength is inherently linked to the proportion of significant betas; more precisely, it estimates the growth rate of significant betas relative to the sample size,  $n$ . The estimator proposed is reported here for convenience, see Chapter 1 for a detailed description:

$$\hat{\alpha} = \begin{cases} 1 + \frac{\ln \hat{\pi}_{nT}}{\ln n}, & \text{if } \hat{\pi}_{nT} > 0, \\ 0, & \text{if } \hat{\pi}_{nT} = 0. \end{cases} \quad (5.2)$$

Where  $\hat{\pi}_{nT}$  is the proportion of statistically significant betas, accounting for multiple testing. This provides an indication of the number of securities effectively influenced by the size factor.

Table 5.2 presents the results of the factor strength estimation over the entire sample period, using the global factor obtained by incorporating the Market-to-Book ratio (MTBV) in the cross-sectional regression phase. Consistent with findings in the literature, the market factor exhibits the highest strength, with a value close to 1. Both the Size and Value factors show substantial strength, with the Size factor slightly outperforming the Value factor. This result suggests that the factors considered are indeed persistent, exerting significant influence across the majority of the securities analyzed.

Table 5.3 displays the strength estimates for the Size factor, evaluated separately for each geographical region. Notably, Europe emerges as the area where the Size factor is most pervasive, with strength exceeding 0.9, indicating a pronounced impact on securities within this region.

<b>Factor</b>	<b>Estimate</b>	<b>Lower Bound</b>	<b>Upper Bound</b>
Market	0.98	0.98	0.98
Size (MV)	0.86	0.86	0.86
Value (MTBV)	0.83	0.83	0.83

TABLE 5.2: Estimated factor strength over the full time period for each of the factors included in the model

<b>Region</b>	<b>Estimate</b>	<b>Lower Bound</b>	<b>Upper Bound</b>
US	0.75	0.75	0.76
EU	0.92	0.91	0.92
BICS	0.79	0.76	0.81

TABLE 5.3: Estimated strength of the size factor (MV) across geographical subsamples

It is reasonable to expect that factor strength may vary over time, potentially leading to periods of strength instability throughout the sample period. To explore this, we estimate strength using 10-year rolling windows with a 6-month step between each consecutive window. This approach provides an overview of how estimated strength evolves over time, potentially capturing shifts in factor pervasiveness across different market periods.

Figure 5.7 displays the estimated strength for each window, along with confidence bands, using the global factor obtained by including the Market-to-Book ratio (MTBV) in the cross-sectional regression phase for the full sample. Results are evaluated both at an aggregate level and a local level, with securities grouped by geographical area for the local analysis.

Figure 5.8 illustrates similar estimates, but here strength is derived from local factors, computed using only geographical subsets of the sample in the cross-sectional regression phase.

The confidence bands are notably narrow across all configurations, indicating a high level of precision in our strength estimates. An exception occurs at the start of the last panel in Figure 5.7, which covers the BICS subset. This wider confidence band is likely due to the smaller sample size for this region in that period, resulting in greater uncertainty.

Comparing Figures 5.5 and 5.6 provides valuable insights. While local factor estimates are somewhat less relevant for our primary objective, they do reveal that different markets exhibit varying exposures and behaviors in relation to the size effect. Notably, the estimated strength for local factors tends to be higher than for the global factor, as local factors are derived from an optimization procedure specific to each regional subsample.

The primary focus of this analysis is the global size factor, illustrated in Figure 5.5. The black line, representing factor strength across the entire set of available securities, suggests that the size factor is semi-strong, with an average value of approximately 0.7. Notably, the most valuable insight here is the clear downward trend, indicating that the influence of the size factor has weakened over time. This decline, combined

with the increasing influence of mega-caps companies, suggests the need to reassess our understanding of this factor, as its dynamics and behavior appear to be evolving.

Further insights emerge when examining geographical subsets. In the European subsample, the factor is particularly pervasive, with strength patterns closely paralleling those observed in the full sample, although the EU exhibits a slightly higher estimated strength than the aggregate. In contrast, the US subsample shows a declining trend in strength from the midpoint of the observation period onward, remaining somewhat below the full sample estimate until a gradual re-alignment near the end of the period. The BICS sub-sample, however, displays the greatest volatility in factor strength, reflecting a higher degree of instability. Furthermore, the BICS sub-sample consistently shows the lowest strength values overall, with estimates remaining below 0.6 throughout the observed period.

In summary, the size factor derived via the cross-sectional approach exhibits a declining trend in strength over time, suggesting a gradual reduction in securities' sensitivity to it. Regionally, the factor exerts its strongest influence within Europe, while the BICS sub-sample shows minimal exposure and heightened instability. The US market displays moderate exposure to the size factor, with a declining trend that eventually converges with full-sample levels. These observations underscore the evolving nature of the size factor and suggest that regional variations play a significant role in its application and impact.

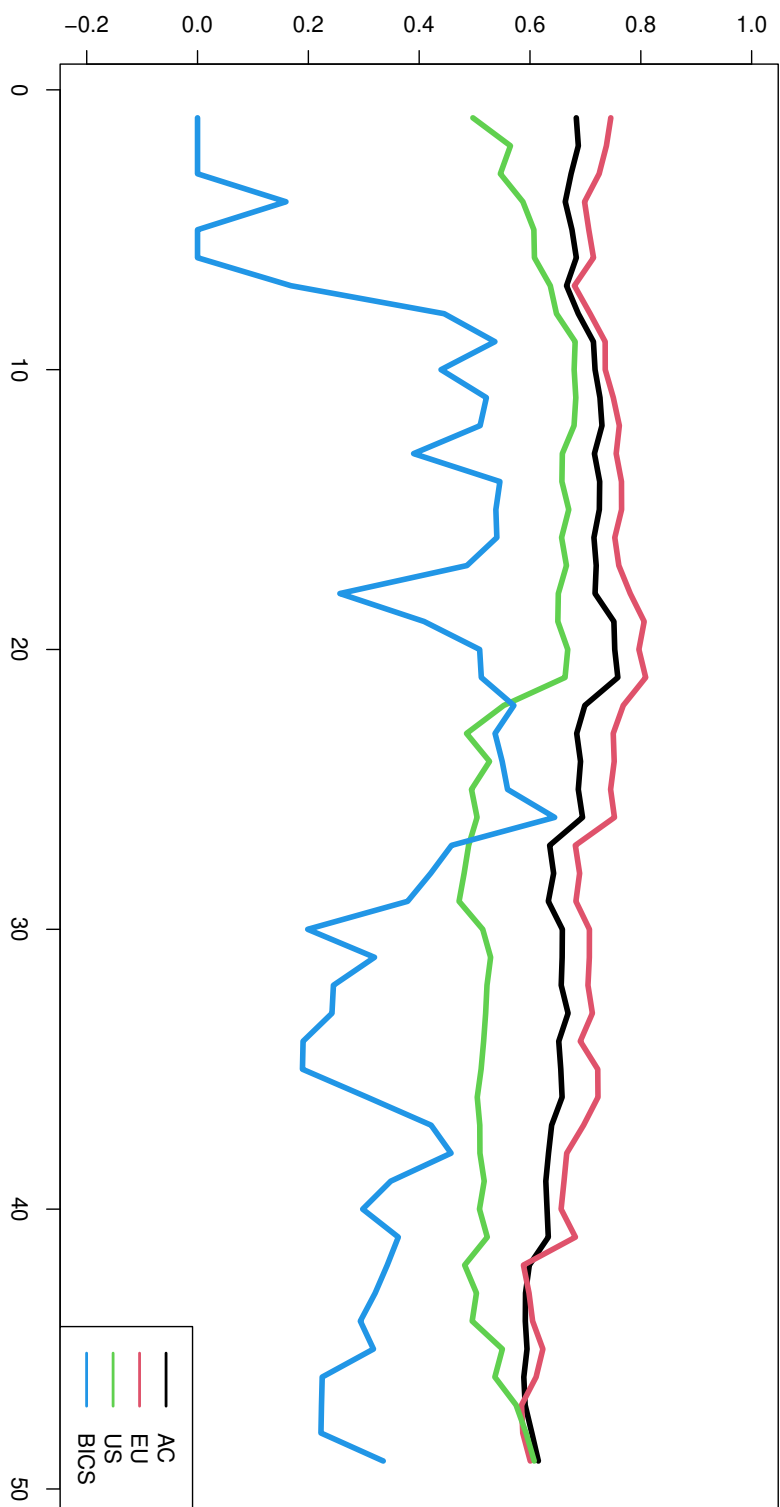


FIGURE 5.5: Estimated strength for the global CS size factor evaluated for each geographical area, 10 years rolling windows with 6m step (Full Sample from 1990).

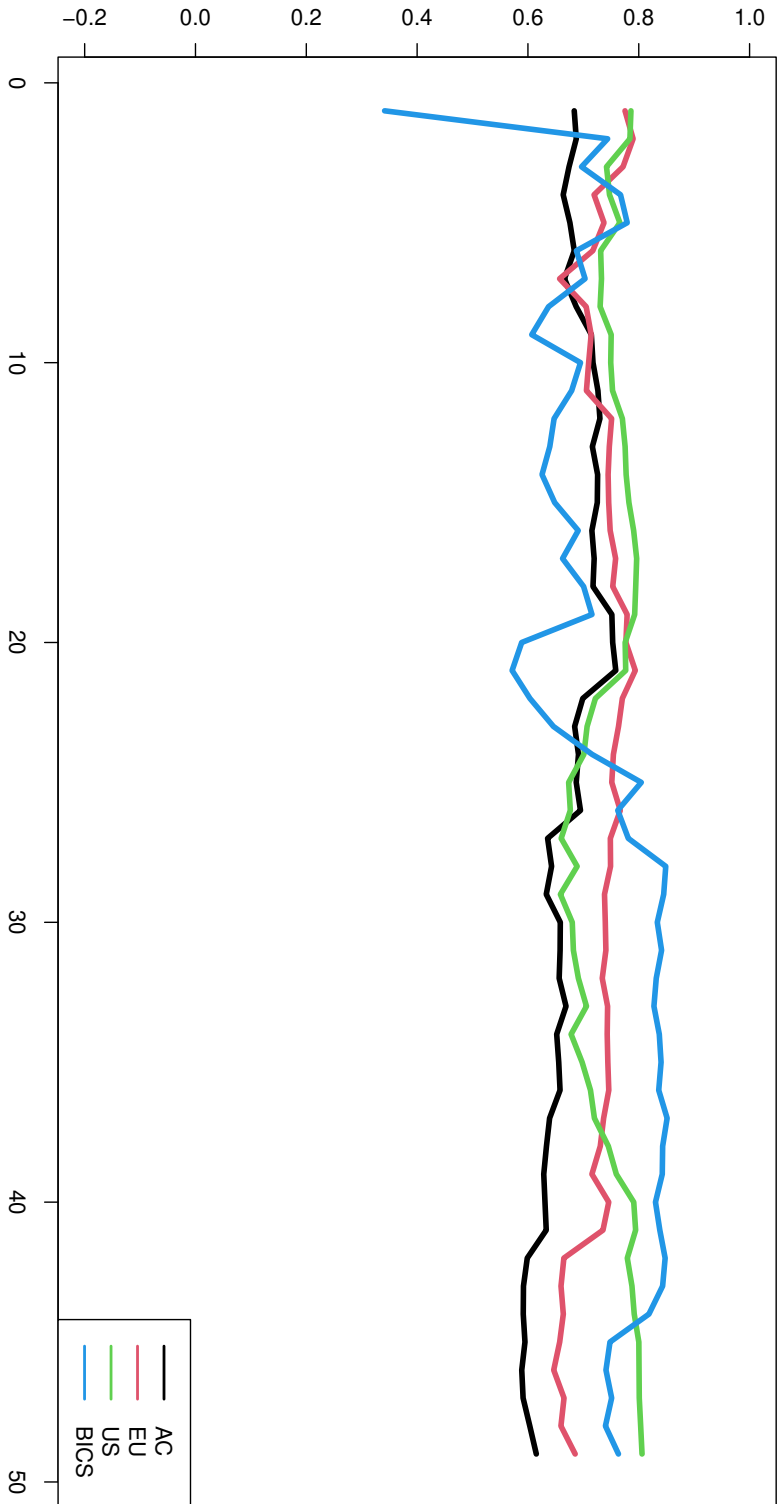


FIGURE 5.6: Estimated strength for the local CS size factor evaluated for each geographical area, 10 years rolling windows with 6m step (Full Sample from 1990).

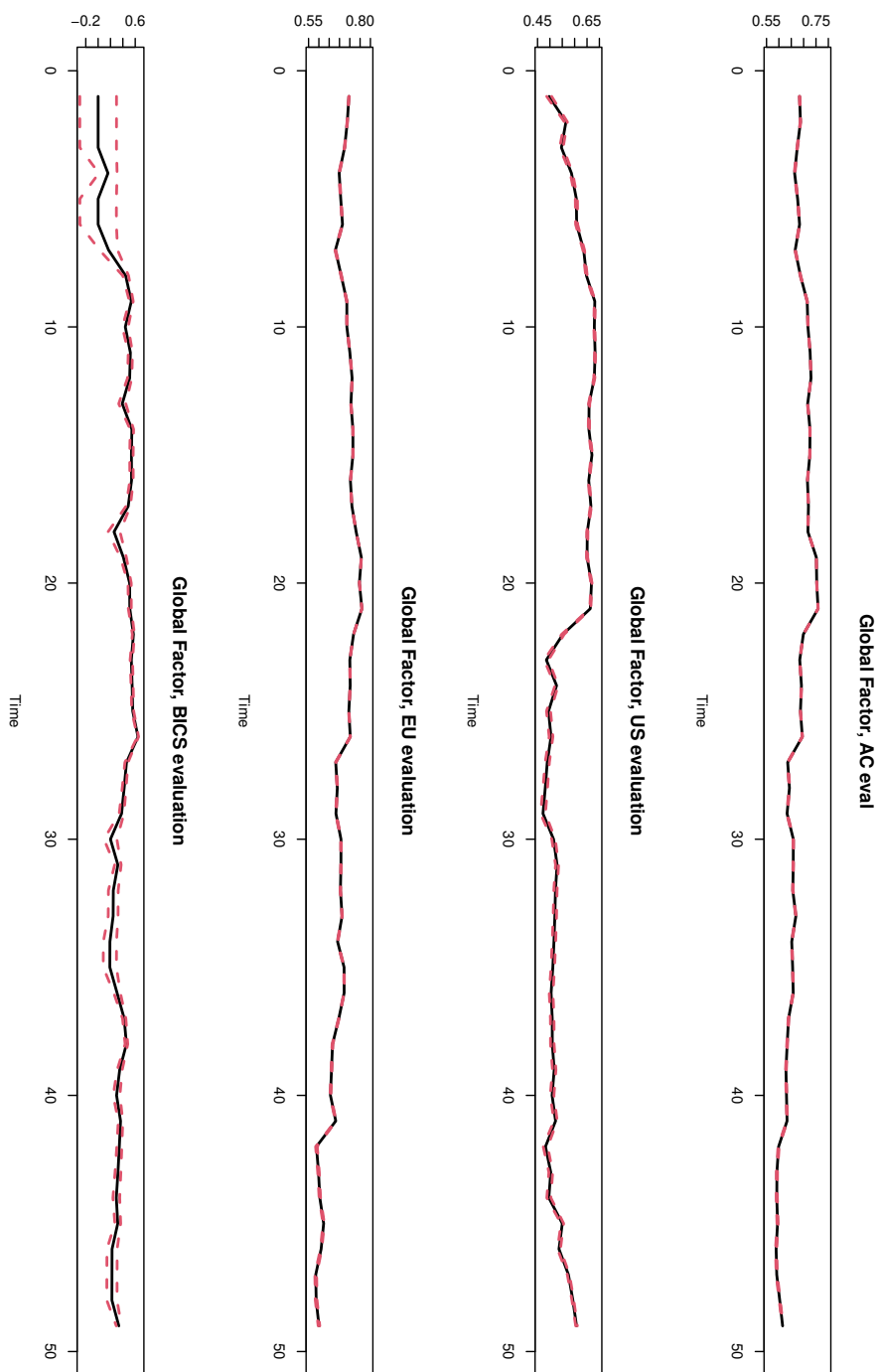


FIGURE 5.7: Estimated strength for the local CS size factor and confidence bands (red dashed line) evaluated for each geographical area, 10 years rolling windows with 6m step (Full Sample from 1990).

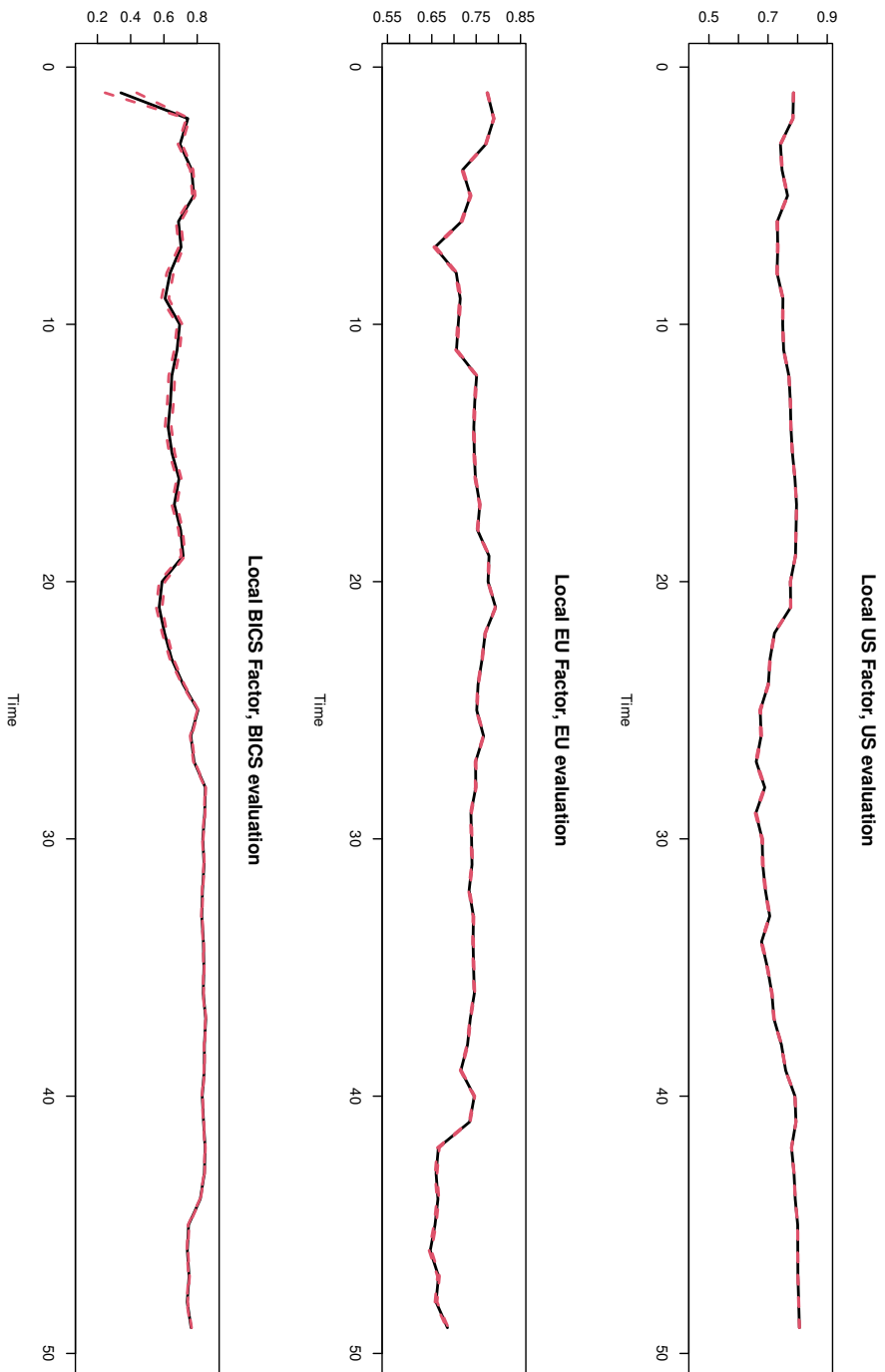


FIGURE 5.8: Estimated strength for the local CS size factor and confidence bands (red dashed line) evaluated for each geographical area, 10 years rolling windows with 6m step (Full Sample from 1990).

# Chapter 6

## Assessing Methodological Weaknesses and Robustness

In this work, we have followed the guidelines dictated by some of the most recent and promising approaches in asset pricing and risk premium estimation. The framework and methodologies employed here are widely adopted, and there is a strong consensus within both academics and professionals that these techniques represent a meaningful advancement beyond traditional, well-established methods that have been intensively utilized in this field.

The aim of this chapter is to explore and examine in greater depth the potential weaknesses of the procedures used, offering a critical analysis of important aspects that are frequently overlooked in the literature. By addressing these limitations, we not only deepen our understanding of the constraints inherent in our approach but also provide insights for future improvements and the development of even more robust methods.

### 6.1 On the robustness of the treatment of outliers

The treatment of outliers is an essential component of any robust analysis, as it helps to ensure that conclusions are not skewed by the presence of abnormal observations. This is especially pertinent when working with financial data, which often contains measurement errors and implausible values that can distort results.

In the literature regarding the methods we have employed, winsorization is a commonly used technique to address this issue, typically involving the removal of 1% of the data. Once the percentage (or, more generally, the number of observations) to be removed is determined, let this be  $x\%$ , the top and bottom  $\frac{x}{2}\%$  of observations with the

highest and lowest values for a chosen variable are removed, respectively. This approach is straightforward, establishes a standardized method for handling outliers, and reduces the distortion they introduce.

However, despite its advantages, this method is not without limitations. Winsorization risks excluding potentially informative observations, which could itself introduce biases into the analysis. This raises uncertainty around determining the optimal percentage of data to exclude, or equivalently, the appropriate threshold values to apply. This question is far from trivial, and it is challenging to identify a “golden rule” that would be universally optimal.

Figure 6.1 illustrates the impact of different winsorization thresholds on the resulting size factor. Here, we examine three thresholds, applying the winsorization scheme to each variable included in the regression, with the exception of Market Value (MV). The sensitivity of results to the chosen threshold is apparent, highlighting how varying thresholds can lead to divergent conclusions. This variability emphasizes the subjective nature of threshold selection and how it can drive the differences in outcomes across analyses.

One potential solution to address this issue is to employ robust regression methods instead of Ordinary Least Squares (OLS) during the cross-sectional phase of the procedure. This choice could lead to significant improvements in the validity of the estimates, as robust regression techniques are designed specifically to mitigate the influence of outliers. By directly addressing the problem of atypical observations, this shift in methodology could also help manage another common issue: the heteroscedasticity of residuals. The OLS framework relies on several assumptions, one of which is homoscedasticity, or constant variance of the residuals. While this is not our primary concern in the current study, as we have not focused on direct inference of the estimated factors, incorporating a robust framework could nonetheless improve the stability and reliability of the results.

Although this approach appears promising, there are both practical and theoretical considerations that must be addressed before asserting its overall benefits to the methodology. On the practical side, robust regression methods are generally more computationally intensive than OLS; however, with advancements in computational power, this is a manageable issue.

From a theoretical perspective, adopting robust regression introduces certain challenges. The procedure employed in this study is built upon desirable properties that ensure consistency with established methods, particularly those proposed by Fama and

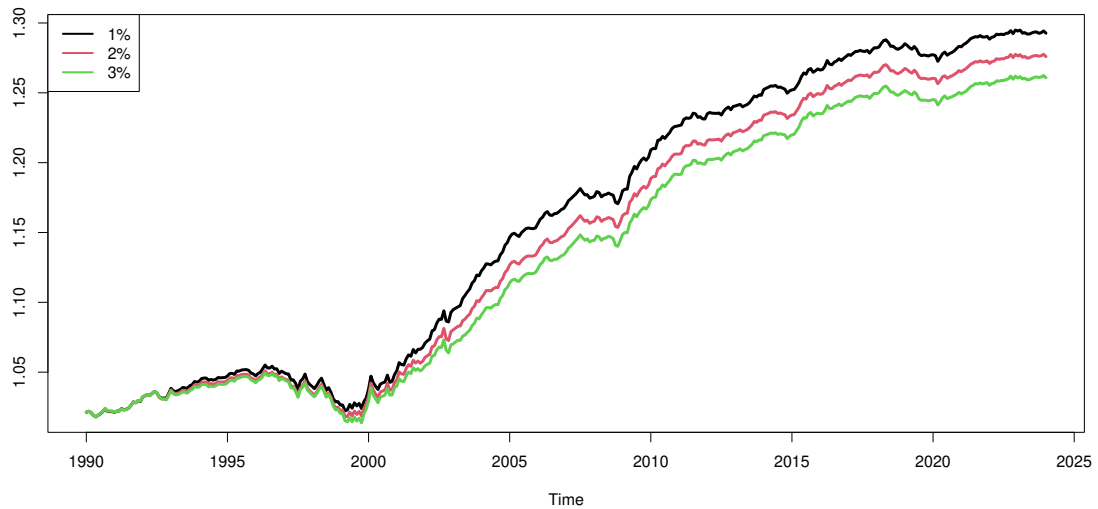


FIGURE 6.1: Comparison of the cumulative size factor obtained using different win-size thresholds, indicated by the total number of excluded observations. The variables included in the model are: MV, MTBV, IG, and OP.

French. Perhaps the most crucial feature of the methodology presented here is its interpretation of regression estimates as portfolio returns, aligning with the conceptual framework of Fama and French's time-series factors. A shift to robust regression might impact this interpretative aspect, potentially altering the fundamental characteristics of the resulting factors. Therefore, further investigation would be required to assess whether this change would preserve the intended coherence of the approach, or if it could lead to substantial differences in factor properties.

## 6.2 On the omitted variable bias

Another issue worth addressing is the omitted variable bias. In Chapter 4, we have already observed how including different explanatory variables in the model specification of cross-sectional regressions affects the resulting size factor. The variation in estimated regression coefficients due to the inclusion of additional variables can be attributed to two main factors.

First, as previously discussed, incorporating more variables in this framework reduces data coverage, decreasing the number of observations available for fitting the regression. Given the properties of the methods used, this reduction changes the composition of the portfolios from which returns are mimicked by the factor of interest.

The second reason, and the primary focus of this section, stems from the omitted variable bias. Omitting relevant variables from a regression model can lead to significant biases in the estimated coefficients. This bias arises when an excluded variable is both a determinant of the dependent variable and correlated with one or more included independent variables, which distorts the coefficient estimates for those included variables.

To illustrate this, consider a simple linear regression model with two independent variables,  $X_1$  and  $X_2$ :

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon \quad (6.1)$$

If  $X_2$  is omitted, the estimated regression becomes:

$$Y = \tilde{\beta}_0 + \tilde{\beta}_1 X_1 + u \quad (6.2)$$

where  $u$  represents the error term, now encompassing the effect of  $X_2$ . The omission of  $X_2$  introduces bias in the estimate of  $\tilde{\beta}_1$ , which can be expressed as:

$$\tilde{\beta}_1 = \beta_1 + \beta_2 \cdot \frac{\text{Cov}(X_1, X_2)}{\text{Var}(X_1)} \quad (6.3)$$

This formula shows that the estimate of  $\tilde{\beta}_1$  will be biased if  $X_1$  and  $X_2$  are correlated and  $\beta_2$  is non-zero. In other words, the effect of  $X_2$  is mistakenly attributed to  $X_1$ , leading to an inaccurate estimation of the coefficient for  $X_1$ .

To mitigate omitted variable bias, it is essential to include all relevant variables that influence the dependent variable and are correlated with the included independent variables. In our case, however, this is particularly challenging: adding more variables reduces the sample size, which in turn alters the portfolio composition. This scenario

presents a trade-off between reducing bias by including more control variables and maintaining adequate data coverage. As a result, we may have to tolerate some degree of distortion in the resulting factor due to this trade-off.

Moreover, we can only control for variables that are actually observed. We lack insight into how unobserved (or latent) variables might influence the factor under examination. This adds another layer of complexity and uncertainty to the validity of the resulting factor. In addition to the impact of omitted observable variables, there may also be latent variables affecting the outcomes.

Some methods to address this issue are suggested in the literature, such as Supervised Principal Component Analysis (SPCA) (Giglio et al., 2023). However, most proposed techniques are employed in time-series settings, with limited or no solutions tailored to the cross-sectional approach used in this study. This limitation highlights an important area for further research, as the biases introduced by omitted variables in cross-sectional models implies the need for more robust solutions.

This issue is particularly important and should always be considered when conducting analyses involving the methods applied here. Acknowledging these limitations is crucial for the accurate interpretation of results and underscores the need for continued research into methods that can better handle omitted variable bias, especially in cross-sectional frameworks.



# Conclusion

This work has presented a comprehensive empirical analysis of the small-cap effect, reassessing its role in asset pricing models and tracking its evolution over recent years. By adopting an alternative approach to factor construction, particularly through cross-sectional regressions, we aimed to produce a more nuanced and “pure” representation of the size factor, isolated, at least partially, from the confounding influences of other variables.

Our findings reveal that while the small-cap effect has traditionally contributed positively to equity returns, its influence has diminished in recent years, aligning with the growing dominance of mega-cap stocks. This shift hints at a possible structural change in the dynamics of the size factor itself. In particular, our results indicate that, in recent years, large-cap firms have increasingly driven market returns, a trend that traditional small-cap strategies may not fully account for.

The cross-sectional approach employed here allowed us to include multiple characteristics in the regression framework, generating factors that are less biased and potentially more reflective of real-world market dynamics. This approach is especially beneficial in managing multivariate interactions, such as correlations between market value and other firm-specific characteristics. By enabling a level of control that traditional Fama-French sorting methods could not achieve, this framework proved advantageous. Additionally, our geographical and sectoral analyses highlighted substantial differences in the small-cap effect’s magnitude and variation across regions, with the size factor showing greater strength in European markets and exhibiting the lowest stability in the BICS region.

However, the methodology adopted in this study is not without its limitations, particularly regarding robustness. One desirable property of the cross-sectional factor approach, compared to time-series factors, also presents a key critique: the potential susceptibility of the retrieved factors to omitted variable bias. Despite efforts to control for available characteristics, limitations in data availability and the possible existence of latent variables may still affect the accuracy of our size factor estimates.

Another issue pertains to the subjectivity in outlier treatment. We adopted the prevailing “state of the art” method of winsorization to mitigate outlier influence. While widely used in current literature, optimal thresholds for winsorization remain ambiguous, and analyses can be sensitive to these cut-offs. This sensitivity suggests a need for further research into refining outlier adjustment methods to enhance the robustness of the methodology and reduce subjectivity.

Future research might also explore alternative approaches, such as robust regressions in the cross-sectional framework, as a means to address outlier influence. Although robust regressions could help mitigate some issues, this would require careful examination due to potential implications for factor properties, particularly the interpretability of factors as long-short portfolio returns, a foundational aspect of our chosen methodology.

In conclusion, this study finds that the small-cap effect remains a relevant but evolving element in asset pricing. Our findings reinforce the importance of periodically re-evaluating factor models to reflect shifting market dynamics. Overall, these results highlight the need for ongoing research into more complex and resilient models capable of capturing the dynamic nature of factors like size, which can vary in impact depending on regional, sectoral, and temporal contexts.

# Appendix

## Proof from Chapter 1

This proof has been recovered from Yang et al. (2024).

For illustration purposes, assume we have a one-factor structure in the cross-section regression, a constant plus one MV property at time  $t-1$ . All MV scores are standardized to z-score (with 0 mean and a standard deviation of 1) among all companies in time  $t$ . In concrete, we have

$$\begin{cases} r_{it} = F_{0t} + F_{MV,t}MV_{i,t-1} + \eta_{it} \\ \mu_{MV,t-1} = \frac{\sum MV_{i,t-1}}{n} = 0 \implies \sum MV_{i,t-1} = 0 \\ \sigma_{MV,t-1}^2 = \frac{\sum (MV_{i,t-1}-0)^2}{n} = 1 \implies \sum (MV_{i,t-1})^2 = n \end{cases} \quad (\text{A.1})$$

The matrix form of equation (A1) is:

$$R_t = X_{t-1}\Gamma_t + \Phi_t, \quad (\text{A.2})$$

where  $R_t$  is the return matrix with a size of  $(n \times 1)$  :  $R_t = [r_{1t}, r_{2t}, r_{3t}, \dots, r_{nt}]^T$ ;  $X_{t-1}$  is the variable matrix with a value of time  $(t-1)$  with a size of  $(n \times 2)$  :

$$X_{t-1} = \begin{bmatrix} 1 & MV_{1,t-1} \\ 1 & MV_{2,t-1} \\ \vdots & \vdots \\ 1 & MV_{n,t-1} \end{bmatrix}_{(n \times 2)}. \quad (\text{A.3})$$

The  $\Gamma_t$  is the matrix of factor returns with a size of  $(2 \times 1)$  :  $\hat{\Gamma}_t = [\hat{F}_{0,t}, \hat{F}_{MV,t}]^T$  and the error term  $\Phi_t = [\eta_{1t}, \eta_{2t}, \dots, \eta_{nt}]^T$ . The Ordinary Least Square (OLS) estimation of  $\Gamma_t$  is:

$$\Gamma_t = (X_{t-1}^T X_{t-1})^{-1} X_{t-1}^T R_t \quad (\text{A.4})$$

$$= \left( \begin{bmatrix} 1 & \cdots & 1 \\ MV_{1,t-1} & \cdots & MV_{n,t-1} \end{bmatrix} \begin{bmatrix} 1 & MV_{1,t-1} \\ \vdots & \vdots \\ 1 & MV_{n,t-1} \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & \cdots & 1 \\ MV_{1,t-1} & \cdots & MV_{n,t-1} \end{bmatrix} R_t \quad (\text{A.5})$$

$$= \left( \begin{bmatrix} n & \sum_{i=1}^n MV_{i,t-1} \\ \sum_{i=1}^n MV_{i,t-1} & \sum_{i=1}^n (MV_{i,t-1})^2 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & \cdots & 1 \\ MV_{1,t-1} & \cdots & MV_{n,t-1} \end{bmatrix} R_t \quad (\text{A.6})$$

$$= \left( \begin{bmatrix} n & 0 \\ 0 & n \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & \cdots & 1 \\ MV_{1,t-1} & \cdots & MV_{n,t-1} \end{bmatrix} R_t \quad (\text{A.7})$$

$$= \frac{1}{n^2} \begin{bmatrix} n & 0 \\ 0 & n \end{bmatrix} \begin{bmatrix} 1 & \cdots & 1 \\ MV_{1,t-1} & \cdots & MV_{n,t-1} \end{bmatrix} R_t \quad (\text{A.8})$$

$$= \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n r_{it} \\ \frac{1}{n} \sum_{i=1}^n MV_{i,t-1} r_{it} \end{bmatrix} = \begin{bmatrix} R_{0,t} \\ R_{MV,t} \end{bmatrix}^T, \quad (\text{A.9})$$

so that we have:

$$\begin{cases} f_{0,t} = \frac{1}{n} r_{1t} + \frac{1}{n} r_{2t} + \cdots + \frac{1}{n} r_{nt} \\ f_{MV,t} = \frac{1}{n} (MV_{1,t-1} r_{1t} + MV_{2,t-1} r_{2t} + \cdots + MV_{n,t-1} r_{nt}) \end{cases} \quad (\text{A.10})$$

Equation (A.10) says that the MV risk factor ( $f_{MV,t}$ ) is a portfolio, with the standardized MV score as weights. Note that the MV score is standardized cross-sectionally to have zero mean, which means that for high-MV companies the  $MV_{i,t-1}$  is positive and for low-MV companies, the  $MV_{i,t-1}$  is negative. In that sense,  $f_{MV,t}$  is the return difference between high- and low-MV companies. It is the 'performance premium': the additional positive/negative return brought by additional MV performance.

## Additional figures from Chapter 4

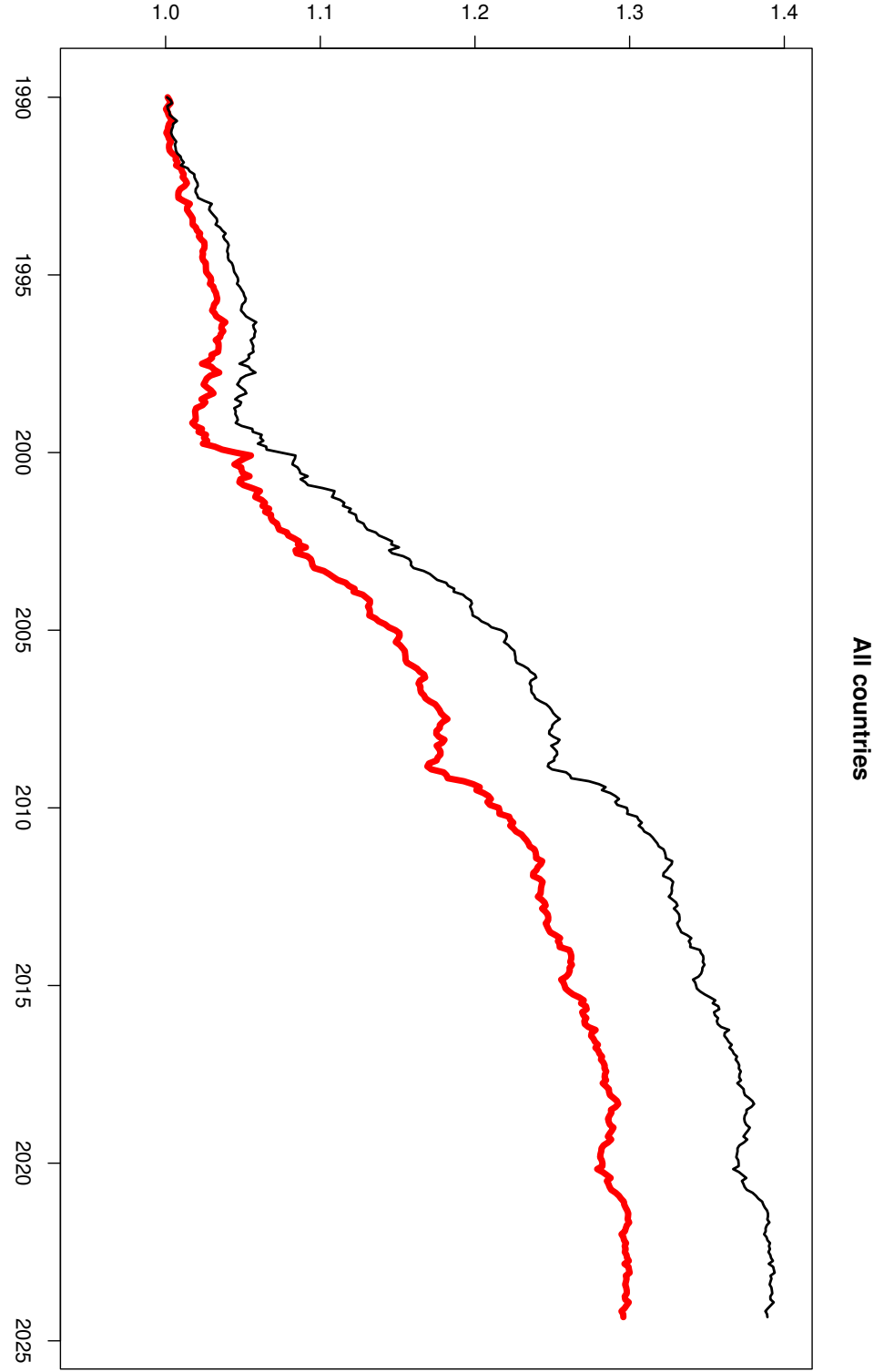


FIGURE .1: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as reference explanatory variable and controlling for the role of Market-to-Book Value (MTBV) - data from 1990 - All countries. Black line: estimate without controlling for MTBV.

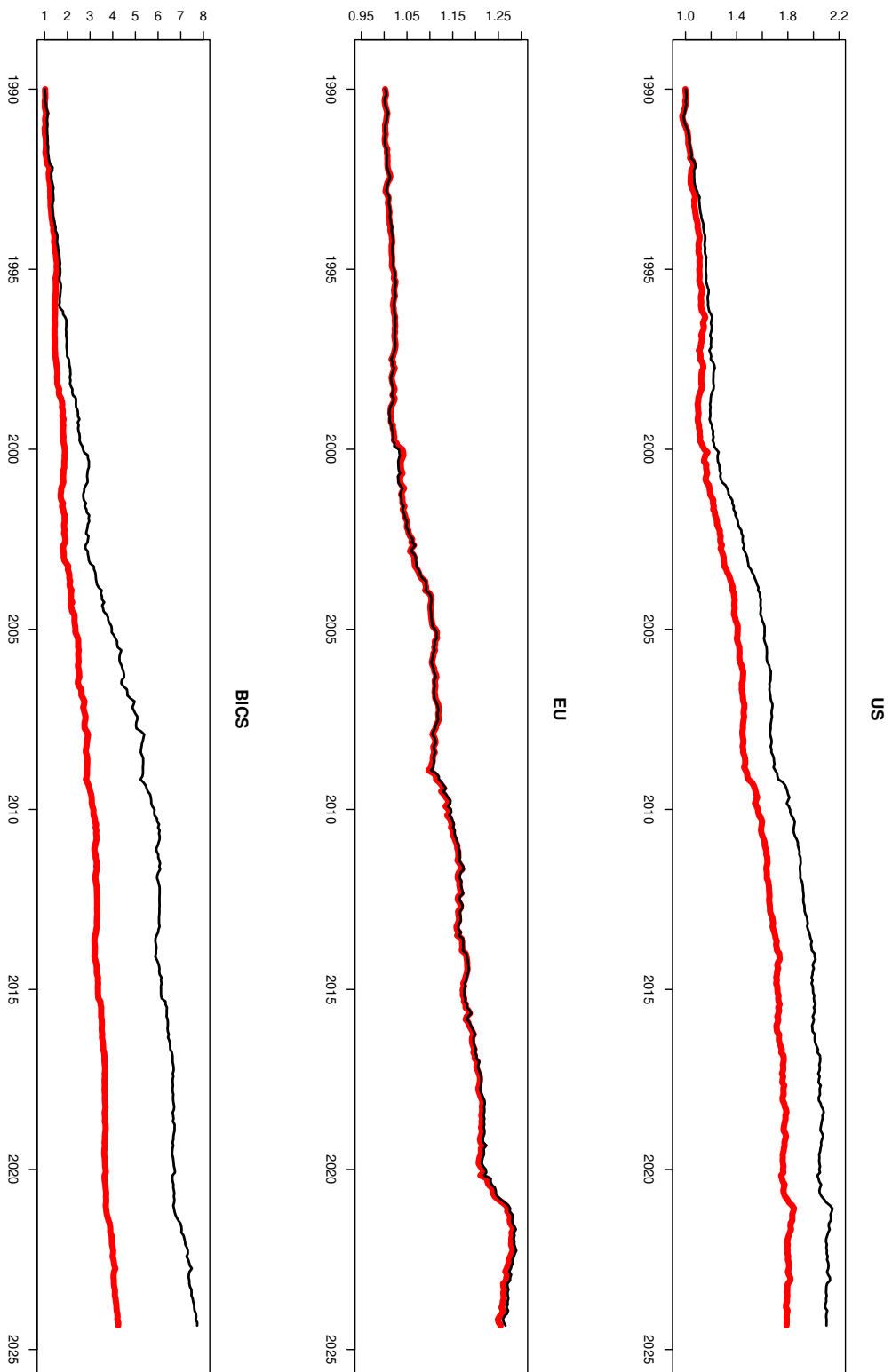


FIGURE 2: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as reference explanatory variable and controlling for the role of Market-to-Book Value (MTBV) - data from 1990 - Estimates made at the geographical group level (i.e., the Size factor is estimated using companies belonging to a specific geographical area). Black line: estimate without controlling for MTBV.

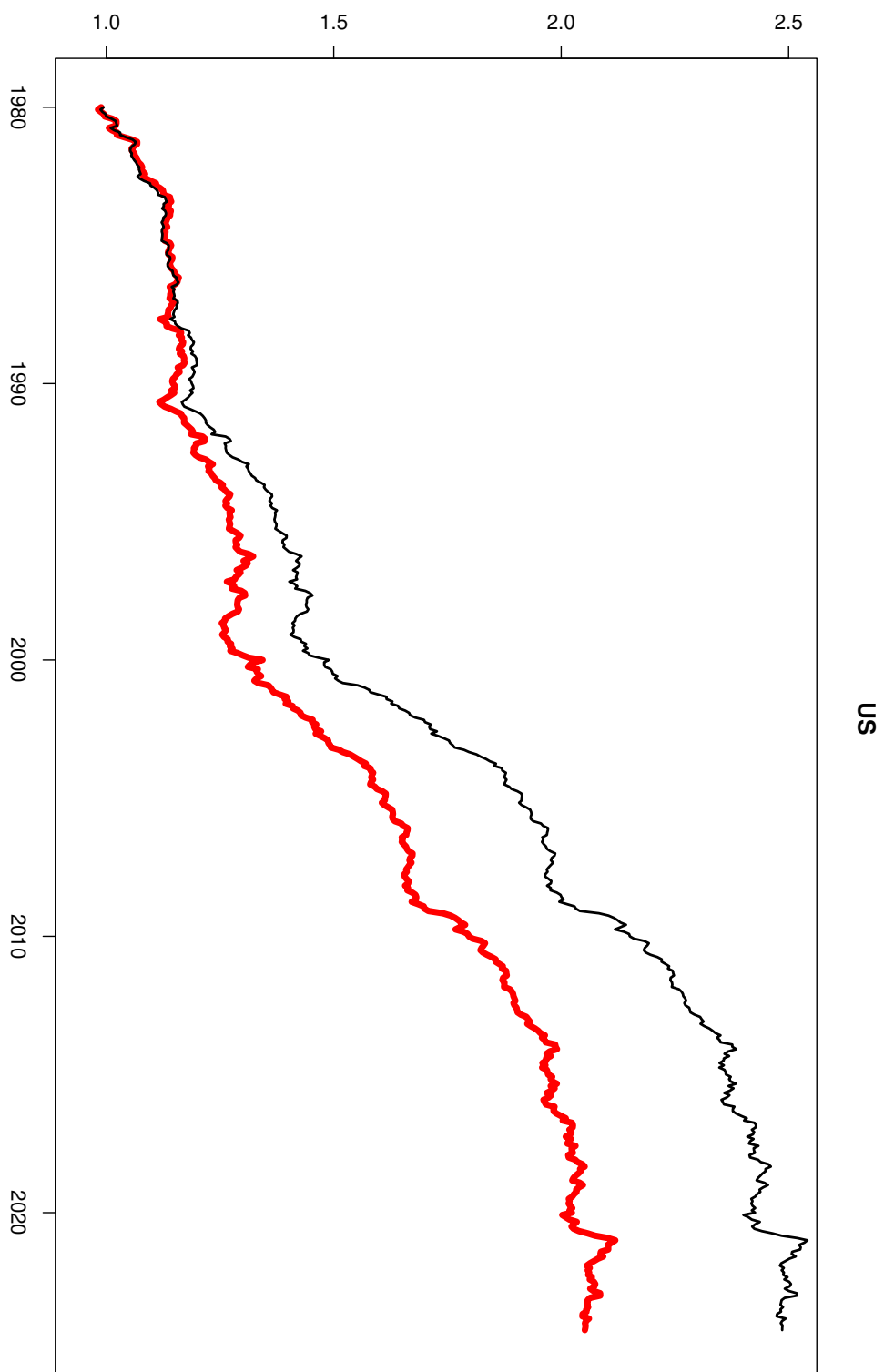


FIGURE .3: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as reference explanatory variable and controlling for the role of Market-to-Book Value (MTBV) - data from 1980 - United States. Black line: estimate without controlling for MTBV.

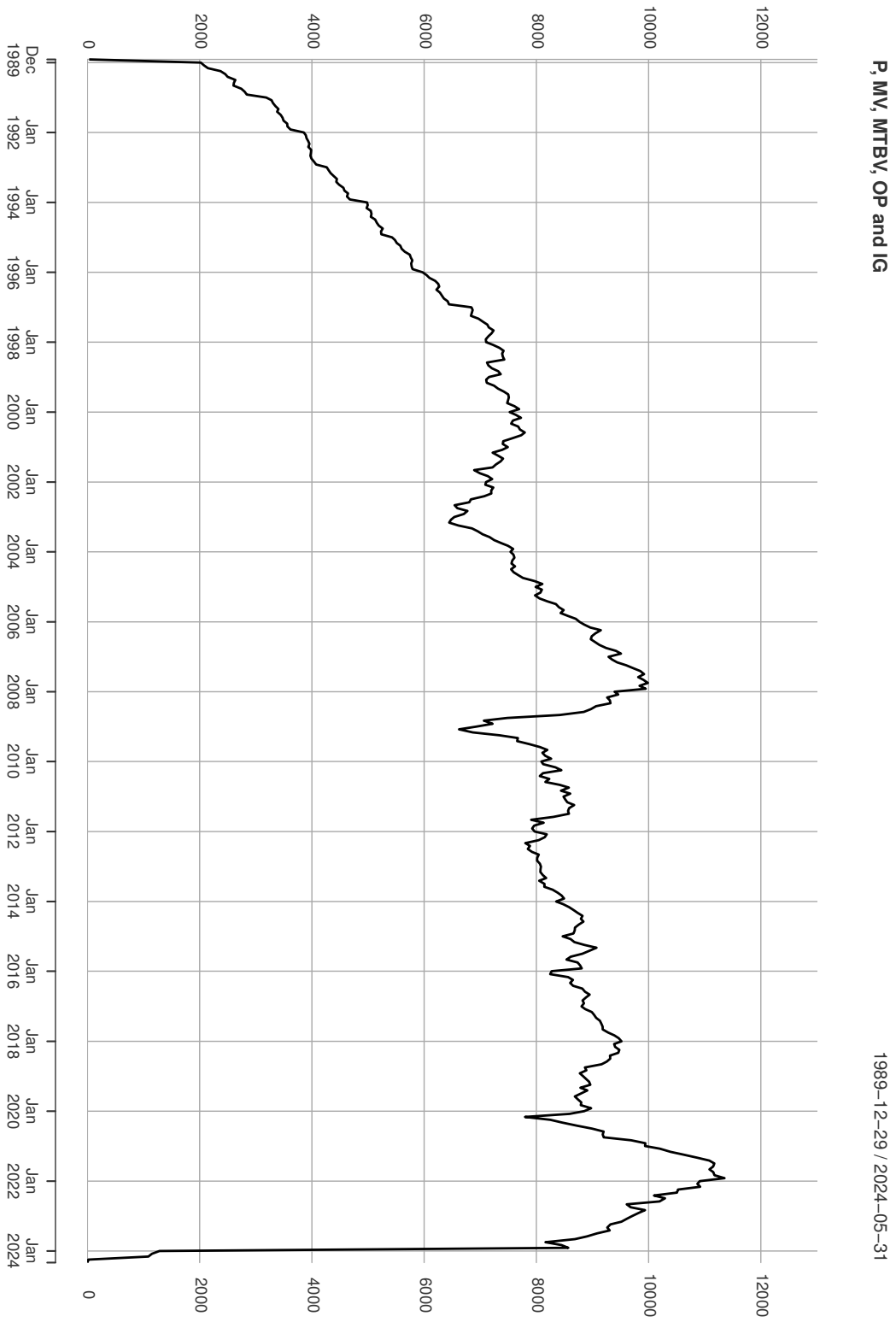


FIGURE .4: Dataset coverage, number of companies for each month, when Price (P), Market Value (MV), Market-to-Book value (MTBV), Operating Profit (OP), and Investment Growth (IG) are all available.

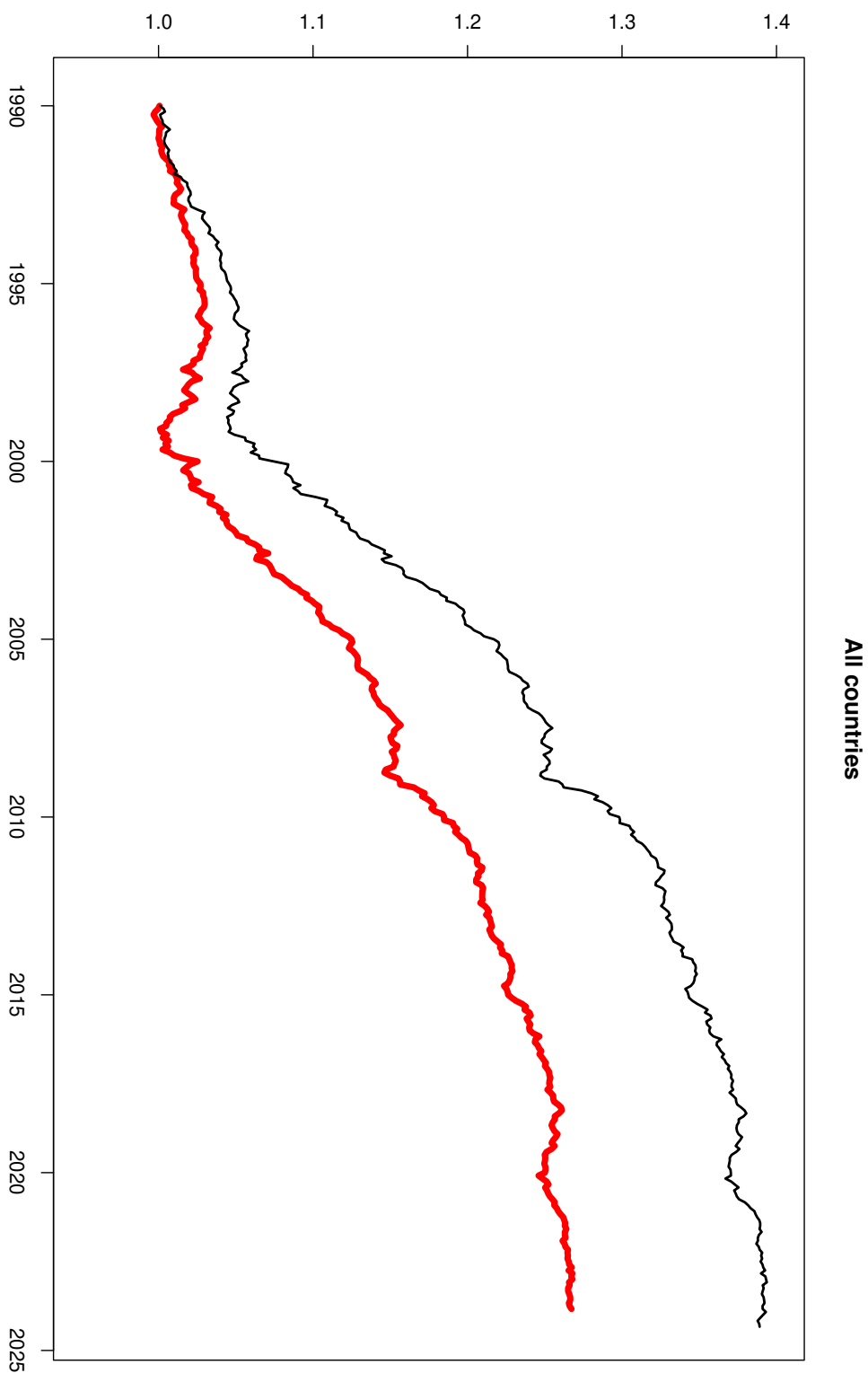


FIGURE .5: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as reference explanatory variable and controlling for the role of Market-to-Book Value (MTBV), Operating Profit (OP), and Investment Growth (IG) - data from 1990 - All countries. Black line: estimate without controlling for other variables.

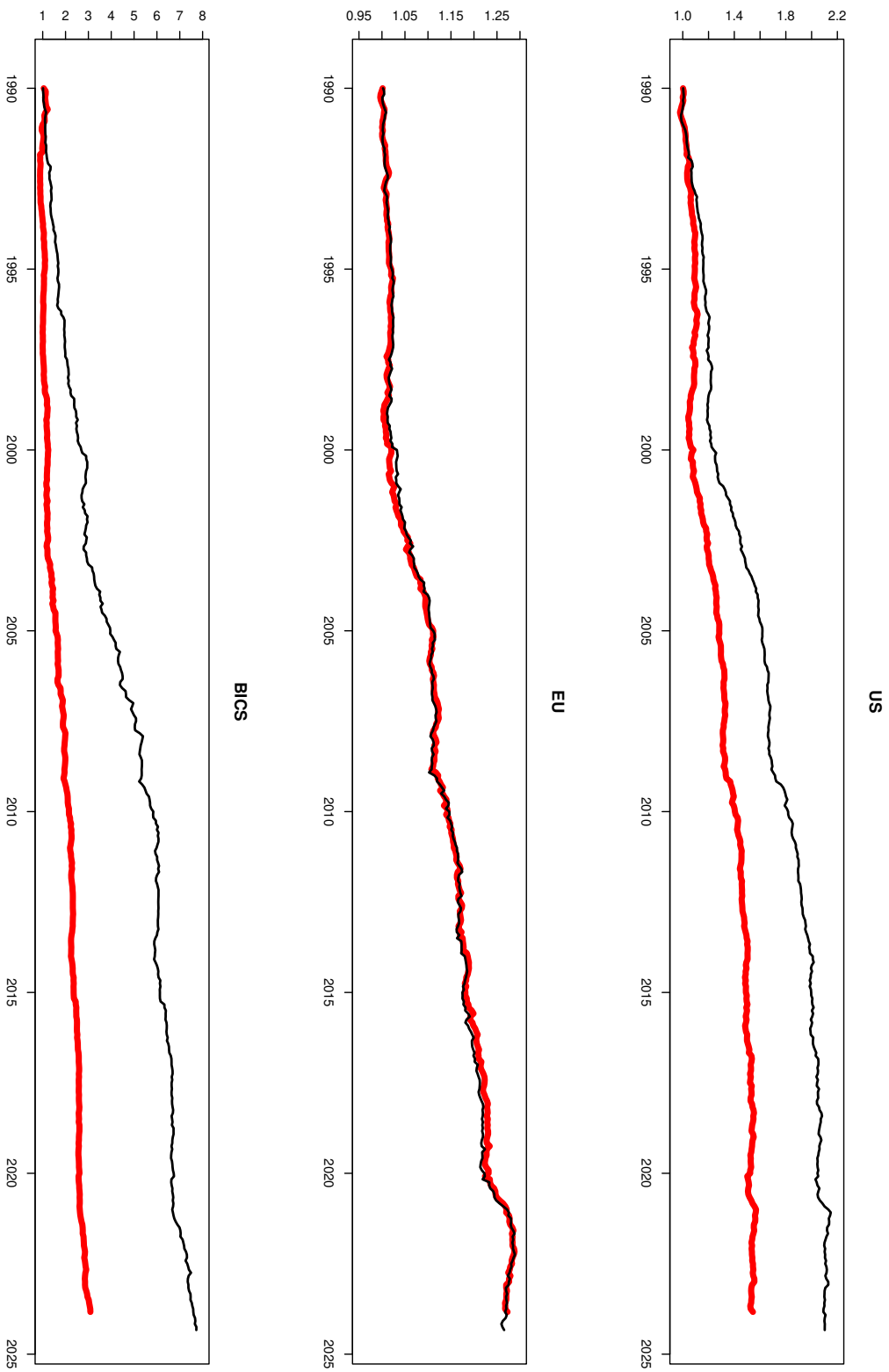


FIGURE 6: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as reference explanatory variable and controlling for the role of Market-to-Book Value (MTBV), Operating Profit (OP), and Investment Growth (IG) - data from 1990 - Estimates made at the geographical group level (i.e., the Size factor is estimated using companies belonging to a specific geographical area). Black line: estimate without controlling for other variables.

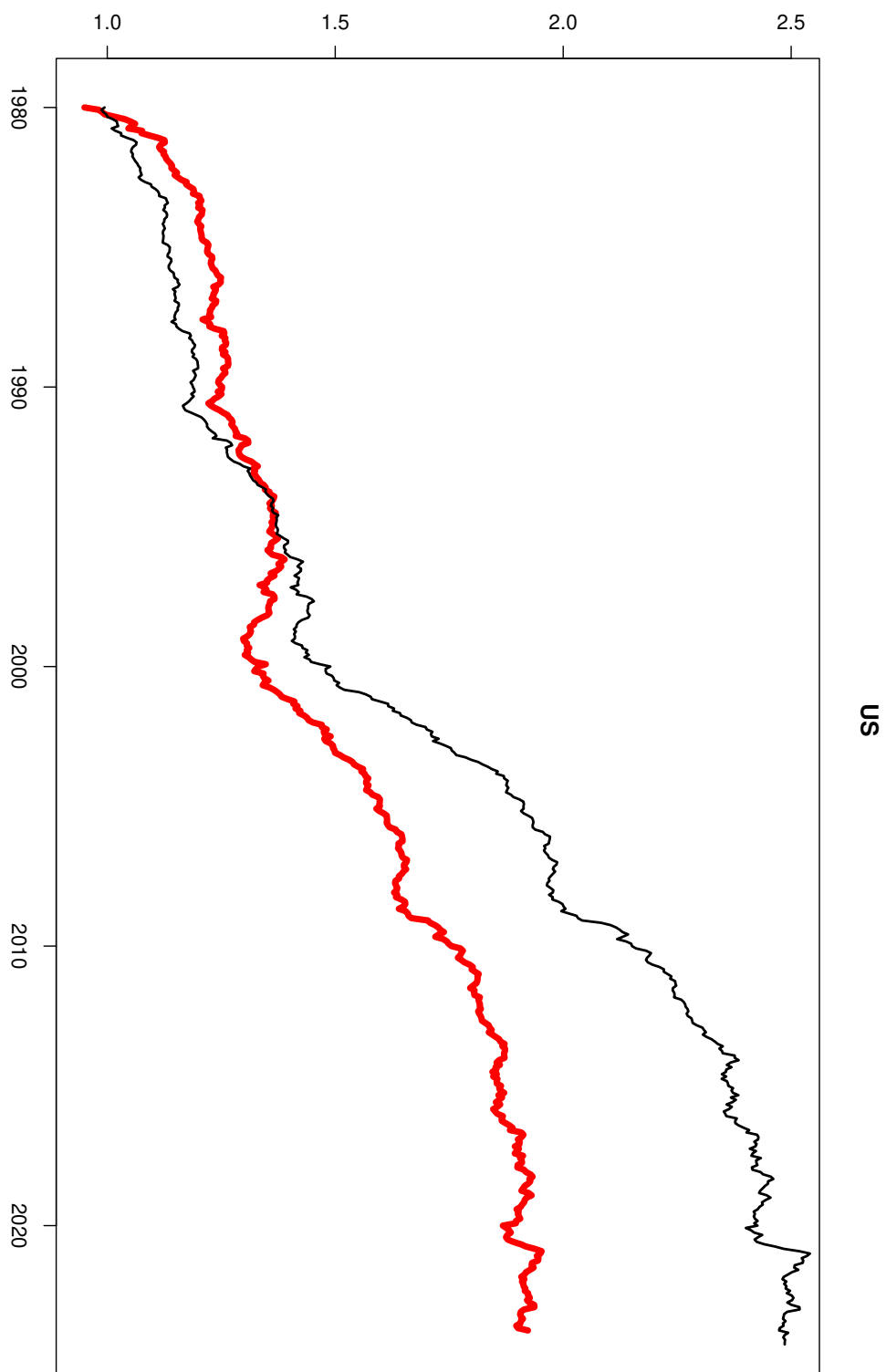


FIGURE .7: Red line: Cumulative returns of the Size factor based on cross-sectional estimates using the Market Value as reference explanatory variable and controlling for the role of Market-to-Book Value (MTBV), Operating Profit (OP), and Investment Growth (IG) - data from 1980 - United States. Black line: estimate without controlling for other variables.



# Bibliography

- Bailey, N., Kapetanios, G., and Pesaran, M. H. (2021). Measurement of factor strength: Theory and practice. *Journal of Applied Econometrics*, 36(5):587–613.
- Fama, E. F. and French, K. R. (1992). The cross-section of expected stock returns. *The Journal of Finance*, 48(2):427–465.
- Fama, E. F. and French, K. R. (2015). A five-factor asset pricing model. *Journal of Financial Economics*, 116(1):1–22.
- Fama, E. F. and French, K. R. (2020). Comparing cross-section and time-series factor models. *The Review of Financial Studies*, 33(5):1891–1926.
- Fama, E. F. and French, K. R. (2023). Production of u.s. rm-rf, smb, and hml in the fama-french data library. Working Paper 23-22, University of Chicago, Booth School of Business.
- Fama, E. F. and MacBeth, J. D. (1973). Risk, return, and equilibrium: Empirical tests. *Journal of Political Economy*, 81(3):607–636.
- Giglio, S., Xiu, D., and Zhang, D. (2023). Test assets and weak factors. Working Paper 21-04, University of Chicago Booth School of Business.
- Lioui, A. and Tarelli, A. (2022). Chasing the esg factor. *Journal of Banking Finance*, 139:106498.
- Pastorello, S. (2001). *Rischio e Rendimento*. Il Mulino.
- Pesaran, M. H. and Smith, R. P. (2021). Factor strengths, pricing errors, and estimation of risk premia. Working Paper 8947, CESifo.
- Yang, R., Caporin, M., and Jiménez-Martín, J.-A. (2024). Esg risk exposure: A tale of two tails. *Quantitative Finance*. Published online, 10 June 2024.

