# ZHAW
# Zurich University of Applied Sciences

School of Management and Law

---

# Bachelor's Thesis

# Estimating Multi-Beta Pricing Models

# With or Without an Intercept:

Further Results from Simulations

---

**Submitted by:** Florin Ernst Akermann

Address: Kleestrasse 3

9472, Grabs

Switzerland

Email address: akermflo@students.zhaw.ch

Student ID: 10-605-517


**Submitted to:** Armin Bänziger

Zurich, 25 May 2016

Zürcher Hochschule
für Angewandte Wissenschaften

**zh aw** School of Management and Law

# Declaration of authorship

I hereby declare that this thesis is my own work, that it has been generated by me without the help of others and that all sources are clearly referenced. I further declare that I will not supply any copies of this thesis to any third parties without written permission by the director of this degree program.

I understand that the Zurich University of Applied Sciences (ZHAW) reserves the right to use plagiarism detection software to make sure that the content of this thesis is completely original. I hereby agree for this thesis, naming me as its author, to be sent abroad for this purpose, where it will be kept in a database to which this university has sole access. I also understand that I will be entitled at any time to ask for my name and any other personal data to be deleted.

Furthermore, I understand that pursuant to § 16 (1) in connection with § 22 (2) of the federal law on universities of applied sciences (FaHG) all rights to this thesis are assigned to the ZHAW, except for the right to be identified as its author.

Student's name (in block letters)

Student's signature

………………………………………..

Zürcher Hochschule
für Angewandte Wissenschaften

**zh aw** School of Management and Law

## Lecturer's statement concerning publication

This statement concerns the publication[1] of the bachelor thesis "Estimating Multi-Beta Pricing Models With or Without an Intercept: Further Results from Simulations".

This bachelor thesis

☐     shall not be published.

☐     shall remain unpublished until (year)      .

☐     is hereby released for unlimited publication.

      ,                            …………………………………...

(Place, date)                             (lecturer's signature)

[1]   "Publication" includes making the thesis available to be read in house, or its distribution to others at cost price or on loan.

## Management Summary

The two-pass method is a common approach for estimating risk premia and examining factor pricing models. It consists of a time-series regression (first-pass) and a cross-sectional regression (second-pass). Two common problems of this approach are the downward bias of the estimates and the large standard error of the estimates. A previous study showed through a simulation approach that the problem of the bias can be mitigated by running at least one of the two regressions without an intercept, whereas the problem of the large standard error can be mitigated by running the second regression without an intercept.

The previous study used a single factor pricing model as the underlying model for its simulation. The objective of this thesis was to provide further evidence for this mitigation method (leaving out the intercepts) by analyzing the mitigating effects in the case of the Fama and French three-factor model. For this purpose, the simulation was based on the simulation of the previous study and was extended to suit the properties of the three-factor model. The simulation comprised two main parts: first, the test data was generated artificially, then the two-pass method was applied on each set of this generated data.

Similar to the findings of the underlying study, it was found that omitting the intercept in at least one of the two regressions decreases the bias of the estimated market premium. Furthermore, omitting the intercept in the cross-sectional regression decreases the standard deviation of the market premium estimates. However, for the two additionally estimated risk premiums (size and value), the mitigating effect on the biases is hardly observable, as the biases of these estimates are already small without omitting the intercept in either of the two regressions. Moreover, the standard errors of the estimates for the size and value premiums did not decrease when the intercept was omitted in at least one of the regressions. In all of the applied variants, the standard error of the estimates for the three premiums are persistently large. Therefore, even with this partially effective mitigation method, it remains difficult to draw statistical conclusions from the two-pass method.

# Table of Contents

# I List of Figures

# II List of Tables

# III List of Codes

# 1   Introduction

At the heart of the asset pricing theory lies the question "why some assets pay higher average returns than others"(Cochrane, 2001, xiii). Within this field of study, linear factor models are the most commonly used models (Cochrane, 2001, p. 229). In fact, the theory originates from the introduction of such a linear factor model (Fama & French, 2004, p. 25). The model in question is the capital asset pricing model (CAPM)[1], introduced by Sharpe (1964) and Lintner (1965). The underlying logic of factor models in asset pricing is that the exposure (beta) to a risk factor should be rewarded with a corresponding premium (gamma)(Cochrane, 2001, p. 81). Based on this logic, factor models are commonly used to estimate risk premiums of stocks and portfolios.

Half a century after its introduction, the CAPM is still one of the most frequently used factor models (Van der Wijst, 2013, p. 81). Since its introduction, it has been extensively scrutinized by empirical researches. Over time, more and more studies showed that the CAPM's ability to explain variations in average returns across assets is limited (Fama & French, 2004, pp. 30-37). Starting with Merton (1973), the idea of extending the CAPM with additional factors emerged (Cochrane, 2001, p. 437). One of the most famous and widely applied "spin-offs" of the scientific discussion about further explanatory factors for the CAPM is the three-factor model introduced by Fama and French (1993, 1996). The big advantage of this newer model is that it includes additional significant risk factors and therefore explains more of the variations of average returns across assets (Cochrane, 2001, p. 437). Throughout the debate about the different factor models, various methods and approaches were used to analyze these models.

A typical method for analyzing linear factor models is to run a cross-sectional regression (Bänziger & Gramespacher, 2015, p. 76). In contrast to a time-series regression, the cross-sectional regression looks at one point or period in time for several assets, whereas the time-series regression is run over time for each asset individually (Cochrane, 2001, pp. 80-81). To put it differently, a time-series regression indicates how much of the variation of the asset return is explained by the explanatory variables (factors) over time[2]. In contrast, the cross-sectional regression indicates how much of the variation of average returns across assets is explained by the risk exposures (betas). These betas are not observable and have to be estimated (Cochrane, 2001, p. 235).

For this purpose, the two-pass method introduced by Black, Jensen, and Scholes (1972) as well as Fama and MacBeth (1973) is often applied (Bänziger & Gramespacher, 2015,

---

[1]In principle, every asset pricing model is a CAPM. However, financial literature refers to the Sharpe-Lintner-Black model as "the" CAPM (Fama & French, 2004, p. 25)

[2]This is the case for the standard formula of time-series regression. Researchers use advanced variants of time-series regressions, which enable to make inferences about average returns between assets (Cochrane, 2001, p. 231).

p. 76). As a first step, the risk exposures of all assets are estimated through a time-series regression (first-pass). In a second step, the derived betas are used as the explanatory variable in a cross-sectional (second-pass) regression in order to estimate the risk premiums (regression coefficient of second-pass regression). Since the explanatory variables of the second-pass regression are measured with error, the estimated coefficients of the second regression are downwards biased (Fama & French, 2004, pp. 30-31). This bias is commonly referred to as the error-in-variables bias. Black et al. (1972) introduced an approach to mitigate this bias, which is now considered a standard approach within empirical factor model studies (Fama & French, 2004, p. 31). They used portfolio returns and estimated portfolio betas instead of single stock betas, since the betas of portfolios can be measured more precisely. However, as simulation studies by Shanken and Zhou (2007, pp. 55-64) have shown, the bias is still quite substantial in the case of small sample sizes. An additional problem of analysing linear factor models with regression analysis are the large standard errors of the resulting estimates (Bänziger & Gramespacher, 2015, p. 76).

## 1.1   Research Question

In a recent study, Bänziger and Gramespacher (2015) addressed these two problems. They tested within a simulation approach whether leaving out the intercepts in at least one of the two regressions mitigates the problem of the downward bias as well as the large standard errors of the estimates. Their motivation was the underlying theory of the linear factor pricing models. This theory suggests that the intercepts of the time-series regression as well as the constant of the cross-sectional regression (Cochrane, 2001, pp. 231-236) should be zero.

The main results of their simulation approach were that leaving out the intercept in at least one of the two regressions decreases the downward bias of the estimates. Moreover, they found that omitting the intercept of the cross-sectional regression also reduces the standard error of the estimates. They conclude that the increase of efficiency and accuracy come at the cost of robustness (Bänziger & Gramespacher, 2015, p. 81).

This thesis aims to provide further evidence on the mitigating effect illustrated by Bänziger and Gramespacher (2015) by applying the two-pass regression method on a three-factor model instead of the CAPM. Therefore, the research question of this thesis is as follows:

"Does, in the case of a multi-factor asset pricing model, leaving out the intercept in at least one of the two regressions decrease the downward bias of the premium estimates and does leaving out the intercept in the second-pass regression decrease the standard error of the premium estimates?"

In order to answer this question, the simulation approach applied by Bänziger and Gramespacher (2015) is extended to the properties of the Fama and French (1993, 1996) three-factor[3] model. Furthermore, since Shanken and Zhou (2007) conducted similar two-pass simulation studies for the three-factor model, but without leaving out the intercepts in the first or second-pass regression, the simulation within this thesis is also considered an extension to their research. These two studies form the main theoretical foundation for the simulation approach of this thesis.

## 1.2   Aims and Focus of the Thesis

Bänziger and Gramespacher carried out their simulation in the statistical programming language $R$. The simulation consists of the data generating part as well as the actual processing of the two-pass method. The aim of this thesis, besides answering the research question, is to provide a comprehensive $R$-code with detailed descriptions (appendix). This aim implies the sub goals of a correct implementation of the data generating simulation, the correct processing of the two-pass method, as well as the embedding of the thesis and its results in the relevant literature[4].

The focus of this thesis and its simulation approach lies solely on indicating and displaying the mitigating effects in the case of a three-factor model. No inferences are made about their statistical significance nor is a hypothesis tested in mathematical terms.

The rest of this thesis is organized in the following way: section 2 describes the most relevant theoretical frame work, which consists of the algebraic description of the two-pass regression method, the error-in-variables problem, the standard approach of mitigating the problem of the large standard errors and error-in-variables, as well as the theoretical background of the multi-factor model to be estimated. Section 3 addresses the methodology applied, namely the parameters calibrated out of the underlying data, as well as the data generating process. In section 4, the results are being presented and discussed. Finally, in section 5, the thesis is summarized and a conclusion is given. In addition, a comprehensive description of the written R-code is provided in the appendix.

---

[3]Any model with three factors could be called the three-factor model. However, similar to the case of the CAPM the financial literature uses "the three-factor model" as synonym for the Fama and French three-factor model. This thesis follows this custom.

[4]The author would like to thank Armin Bänziger for granting him access to their R-code.

## 2    Literature Review

In a first step, the theoretical foundation of the two-pass regression and the *error-in-variables-bias* problem comprised therein are introduced. Thereafter, a common approach to mitigate the problem of the large standard errors and bias of the estimates is illustrated. Finally, the multivariate factor model applied is presented.

### 2.1    Two-pass regression

The two-pass regression method introduced by Black et al. (1972) as well as Fama and MacBeth (1973) has become a standard approach for empirical tests (Fama & French, 2004, p. 31). Since its introductions, various researchers have reviewed and applied this approach. Within this thesis, the description of the two-pass is based on the elaborations of Bänziger and Gramespacher (2015, p. 77) as well as Shanken and Zhou (2007, p. 43-44).

In the first-pass regression, all factor loadings (betas) of the $j$ factors $f$ for $N$ assets are estimated by applying the ordinary least square (OLS) method on equation (1). It describes a multivariate time-series regression model with K regressor.

$$R_{i,t} = \alpha_i + \beta_{i1}f_{1t} + ... + \beta_{iK}f_{Kt} + \epsilon_{it}, \qquad i = 1, ..., N, \qquad t = 1, ..., T, \qquad (1)$$

where

$R_{i,t}$ = return on asset $i$ in period $t$,

$f_{jt}$ = the value of factor $j$ in period $t$ ,    $j = 1, ..., K,$

$\beta_{i,j}$ = the factor loading of asset i on the $j$th factor,

$\epsilon_{it}$ = the random errors, $E(\epsilon_{i,t}) = 0, E(\epsilon_{i,s}, \epsilon_{j,t}) = \sigma_{ij}$ for $s = t$ and $0$ otherwise[5],

$T$ = the number of of time-series observations,

$\alpha_i$ = pricing error (Cochrane, 2001, pp. 230-231).

As shown by Wooldridge (2013, p. 802), the estimation of the $K$ coefficients for $N$ assets by applying OLS to equation (1) can be written as

$$\hat{\beta} = (F'F)^{-1}F'R, \qquad (2)$$

where

$R = T \times N$ matrix containing $T$ returns for $N$ assets each,

---

[5]As is shown in section 3.3, it is assumed that the error terms are independent over time but not across assets (Bänziger & Gramespacher, 2015, p. 77)

$F = T \times K$ matrix containing $T$ sets of $K$ factor values.

Equation (2) results in a $K \times N$ matrix of beta estimates ($\hat{\beta}$). For the second-pass, model (3) is estimated by running a second OLS on the regression model (5). The general expected return equation of a linear regression model is

$$E(R_i) = \gamma + \gamma_1 \beta_1 + ... + \gamma_K \beta_{iK}, \tag{3}$$

where

$\gamma_j$ = risk premium of factor $f_j$,

$\gamma$ = common constant for all assets,

all other variables defined as in equation (1).

Bänziger and Gramespacher (2015, p. 77) show that through the assumption that a risk free rate $(r_f)$ exists, regression model (3) is adjusted to model (4). Because $\gamma$ has no beta-exposure, it must be risk-free and therefore equal to $r_f$. Thus, model (3) can be rearranged as follows:

$$E(R_i) = \gamma + \gamma_1 \beta_1 + ... + \gamma_K \beta_{iK},$$
$$E(R_i) - r_f = E(R_i^e) = \gamma_1 \beta_1 + ... + \gamma_K \beta_{iK}, \tag{4}$$

where $E(R_i^e)$ is the expected excess return of asset $i$. The estimated betas $\hat{\beta}$ are used as the independent variables for the second multivariate OLS applied on equation (5) in order to derive an estimator of $\Gamma = (\gamma_0, \gamma_1, ..., \gamma_K)'$.

$$\bar{R}_i^e = \gamma_0 + \gamma_1 \hat{\beta}_{i1} + ... + \gamma_K \hat{\beta}_{iK} + u_i, \tag{5}$$

where

$\bar{R}_i^e = \frac{1}{T} \sum_{t=1}^{T} R_{it}^e,$

$\hat{\beta}_{ij}$ = the $(i, j)^{th}$ element of matrix $\hat{\beta}$,

$u_i$ = a random error term,

all other variables are defined as in equation (1).

Two important implications can be drawn from the two-pass method described above: Shanken and Zhou (2007, p. 43) point out that equation (4) implies that the constant ($\gamma$) is zero. Cochrane (2001, p. 230) further states that the comparison of model (1) and model (4) implies that $\alpha$ should be zero in the time-series regression. These two implications are the $H_0$ hypothesis to be imposed within the simulation approach of this thesis. The $H_0$ for

the first-pass regression is $\alpha = 0$. Consequently, the $H_0$ of the second-pass regression is $\gamma_0 = 0$.

As before, the estimation of the coefficients $\Gamma$ can also be derived by using matrix algebra. Following Cochrane (2001, p. 236) as well as Shanken and Zhou (2007, p. 43), the OLS estimation of $\Gamma$ in equation (5) can be described as

$$\hat{\Gamma}_t = (\hat{\beta}'\hat{\beta})^{-1}\hat{\beta}'R^e, \tag{6}$$

where $R^e$ denotes a $T \times N$ matrix containing $T$ excess returns for $N$ assets. The resulting $\hat{\Gamma}_t$ contains $K$ vectors with $T$ estimates of $\gamma_j$ each. Rather than running an OLS for $T$ Vectors with $N$ returns, $R^e$ is replaced by $\bar{R}^e$ ($N$-Vector with means of $T$ sample returns). Therefore, equation (7) results in $K$ estimates ($\hat{\Gamma} = (\gamma_0, \gamma_1, ..., \gamma_K)$) for all $K$ coefficients of the second regression.

$$\hat{\Gamma} = \sum_{t=0}^{T} \hat{\Gamma}_t/T = (\hat{\beta}'\hat{\beta})^{-1}\hat{\beta}'\bar{R}^e. \tag{7}$$

Thus, as can be seen in equations (6) and (7), estimates are being used as explanatory variables for the second-pass regression.

In short, the two-pass method consists of a first-pass (time series regression) and a second-pass (cross-sectional) regression. The theory for both regression models state that the intercept should be zero. These theoretical restrictions pose the basis for the simulation approach within this thesis. Furthermore, it has been shown that the second-pass regression is carried out with estimates as explanatory variables. The implication of this circumstance will be discussed in the subsequent section.

## 2.2   Error-in-Variables

As the explanatory variables ($\hat{\beta}$) of the second linear regression are measured with error, the coefficients of the second regression ($\hat{\Gamma}$) are biased (Brooks, 2014, pp. 236-237). While the direction and magnitude of this bias is determinable through a few steps in the case of a single factor model, attempting the same in the case of a multivariate regression with multiple explanatory variables measured with error is a complicated endeavour (Kennedy, 2008, pp. 167-168). Following Wooldridge (2013, pp. 311-312), an algebraic explanation for the directional bias in the case of a single factor regression model is introduced. His explanations are similar to what Black et al. (1972) wrote originally, but somewhat more detailed. However, the notation of Black et al. (1972) is used as the subsequent section follows their notation.

Therefore, for the sake of consistency, the following example uses $\hat{\beta}$ as the explanatory variable measured with an error. In most of the literature on measurement errors in an explanatory variable, $\beta$ is used as the coefficient of the regression model rather than as the explanatory variable. Moreover, most literature denotes the true unobservable value with a $*$, whereas Black et al. (1972) simply denote the true value of beta as $\beta$ and the beta measured with error as $\hat{\beta}$.

In the case of a single factor model ($K = 1$), the second-pass regression model (5) simplifies to

$$R^e = \gamma_0 + \gamma_1 \beta_1 + u. \tag{8}$$

Index $i$ is omitted in the elaborations within this section. All equations from (8) to (15) address the case of a single asset. Instead of the true measure of exposure to risk $\beta_1$, $\hat{\beta}_1$ is plugged into the second-pass regression model (5). $\hat{\beta}_1$ is $\beta_1$ measured with an error and defined in the following way:

$$\hat{\beta}_1 = \beta_1 - e_1, \tag{9}$$

where $e_1$ is the measurement error of $\beta_1$. Plugging equation (9) into equation (8), it follows that

$$R^e = \gamma_0 + \gamma_1 \hat{\beta}_1 + (u - \gamma_1 e_1). \tag{10}$$

Thus, equation (10) describes the actual properties of the composite error term ($u - \gamma_1 e_1$) of the second-pass regression when the explanatory variable ($\hat{\beta}$) contains measurement error $e_1$. Wooldridge (2013) assumes that the measurement errors ($e_1$) are on average zero across the whole population ($E(e_1) = 0$) and independent of $\beta_1$. This can be stated as

$$Cov(\beta_1, e_1) = 0. \tag{11}$$

Accordingly, Brooks (2014, p. 237) describes $e_1$ as additional noise that is independent of $\beta_1$. Moreover, Wooldridge (2013, p. 320) assumes that $u$ is uncorrelated with $\beta_1$ and $\hat{\beta}_1$. If assumption (11) is true and $e_1 = \beta_1 - \hat{\beta}_1$ then $e$ is correlated with $\hat{\beta}$. Wooldridge shows that

$$Cov(\hat{\beta}_1, e_1) = E(\hat{\beta}_1, e_1) = E(\beta_1, e_1) + E(e_1^2) = 0 + \sigma_{e_1}^2 = \sigma_{e_1}^2. \tag{12}$$

Relation (12) suggests that the variance of the measurement error ($e_1$) is equal to the covariance between $\hat{\beta}_1$ and $e_1$ if assumption (11) holds. Since $u$ and $\beta_1$ are uncorrelated,

the covariance between $\hat{\beta}_1$ and the composite error term $(u - \gamma_1 e_1)$ can be described as follows:

$$Cov(\hat{\beta}_1, u - \gamma_1 e_1) = -\gamma_1 Cov(\hat{\beta}_1, e_1) = -\gamma_1 \sigma_{e_1}^2. \tag{13}$$

In order to determine the magnitude and direction of the error-in-variables bias for a simple linear regression, Wooldridge (2013) introduces equation (14), which is the final formula in the proof that $\hat{\gamma}_1$ is an unbiased and consistent estimator for $\gamma_1$ since $Cov(\beta_1, u) = 0$ (Wooldridge, 2013, p. 173).

$$plim \quad \hat{\gamma}_1 = \gamma_1 + \frac{Cov(\beta_1, u)}{Var(\beta_1)}. \tag{14}$$

However, in the case of the two-pass regression, the aggregated error term of equation (10) is $u_i - \gamma e_1$ rather than solely $u_i$. Therefore, in the case of errors in the explanatory variables, formula (14) is described in the following way by utilizing equation (12) and the fact[6] that $Var(\hat{\beta}_1) = Var(\beta_1) + Var(e_1)$.

$$
\begin{aligned}
plim \quad \hat{\gamma}_1 &= \gamma_1 + \frac{Cov(\hat{\beta}_1, u_i - \gamma_i e_1)}{Var(\hat{\beta}_1)} \\
&= \gamma_1 - \frac{\gamma_1 \sigma_{e_1}^2}{\sigma_{\beta_1}^2 + \sigma_{e_1}^2} \\
&= \gamma_1 \left( 1 - \frac{\sigma_{e_1}^2}{\sigma_{\beta_1}^2 + \sigma_{e_1}^2} \right) \\
&= \gamma_1 \left( \frac{\sigma_{\beta_1}^2 + \sigma_{e_1}^2}{\sigma_{\beta_1}^2 + \sigma_{e_1}^2} - \frac{\sigma_{e_1}^2}{\sigma_{\beta_1}^2 + \sigma_{e_1}^2} \right) \\
&= \gamma_1 \left( \frac{\sigma_{\beta_1}^2}{\sigma_{\beta_1}^2 + \sigma_{e_1}^2} \right).
\end{aligned} \tag{15}
$$

Therefore, since the term in the brackets of the final line in equation (15) is always smaller than 1 and decreases the larger $e$ becomes, $\hat{\gamma}_1$ is more biased the larger the measurement errors ($e$). Moreover, it becomes evident that the estimation is biased towards zero. To put it differently, if $\gamma_1$ is positive, the bias is negative and vice versa (Brooks, 2014, p. 237). Additionally, the "bracket term" shows that the relative size of $\sigma_{e_1}$ to $\sigma_{\beta_1}$ is also of relevance. Ceteris paribus, the larger $\sigma_{\beta_1}$ the smaller the disturbance through $e$ within the simple linear regression.

As has been shown, the magnitude and direction of the error-in-variables problem is

---

[6]This is true due to the assumption introduced in equation (11).

discernible within a bivariate linear regression. However, in the case of multivariate regression, the derivation of these dimensions is a complicated endeavour (Kennedy, 2008, p. 167). Wooldridge (2013, p. 323) points out that assumption (11) does not hold in the case of a multivariate estimation equation and therefore makes the derivation of implications of measurement errors very complex. Thus, the derivation of the magnitude and direction of the error-in-variables bias within a multivariate regression equation with multiple explanatory variables is beyond the scope of this thesis.

In short, this section has shown the fundamental logic of the error-in-variables problem and draws attention to the question of how to mitigate the resulting bias. Accordingly, the relevant theoretical background about a common error-in-variables mitigation method will be outlined in the subsequent section.

## 2.3   A Common Mitigation Method

Since, the focus of this thesis is to extend the simulation by Bänziger and Gramespacher (2015) in one specific way (multi-factor instead of single-factor model), only approaches applied by them will be discussed. Bänziger and Gramespacher made use of the method introduced by Black et al. (1972). It is a very common approach (Fama & French, 2004, p. 31). The method addresses both, the error-in-variables problem as well as the large standard deviation of the estimates (Black et al., 1972, pp. 20-21). This section is structured in the following way: first, the underlying concepts of the method by Black et al. are outlined. Second, the data snooping bias discussed by Lo and MacKinlay (1990) is briefly considered, as it addresses the danger of structuring test statistics in a specific way.

The approach comprises two crucial concepts: first, similar to the findings of Blume (1970), they show that the portfolio beta measurement error is smaller than a single stock beta measurement error. Therefore, the error-in-variables problem can be mitigated by using portfolios rather than single stocks. Second, they show that the sampling variability of the estimates in the second-pass regression is reduced when the betas of the various portfolios as the explanatory variables are well dispersed in the test statistics. To put it differently, using stock groups where the individual stock is assigned to a group based on the stock's estimated risk exposure, yields consistent estimates for the risk premium within the two-pass approach (fixed time-series length, increasing number of assets $N$) (Black et al., 1972, p. 20). This is not the case if the grouping procedure is not applied (Shanken, 1992, p. 3-4). In the following the just described will be outlined in detail.

Black et al. (1972) group $N$ assets into $M$ equally sized portfolios $P$. As they empirically test the case of a single factor model (CAPM), the time-series (first-pass) regression equation (1) shortens to:

$$R_{Pt} = \alpha_P + \beta_{P1}f_{1t} + \epsilon_{Pt}, \qquad i = 1, ..., N, \quad t = 1, ..., T, \quad P = 1, ..., M, \qquad (16)$$

where

$R_{Pit}$ = return on asset $i$ in period $t$ within portfolio $P$,

$f_{jt}$ = the value of factor $j$ in period $t$ (CAPM: $f_t$ = market return of period $t$),

$L = \frac{N}{M}$ = number of assets within each portfolio $P$,

$\epsilon_{Pt} = \frac{1}{L}\sum_{i=1}^{L}\epsilon_{Pit}$ = error term of the portfolio $P$,

$R_{Pt} = \frac{1}{L}\sum_{i=1}^{L}R_{Pit}$.

$R_{Pt}$ denotes the return of portfolio P in period t and is used as the response variable for the OLS regression applied on (16), which estimates $\hat{\beta}_{Pit}$. At this point it should be mentioned that a slightly different notation to Black et al. (1972) is used. They denote $\epsilon$ as $e$, however, since some of the formulas from section 2.2 are reused, $e$ is already assigned to the measurement error of $\beta$ in this thesis.

They assume that $\epsilon_{it}$ are independently distributed and that the variance of $\epsilon$ is constant for all $i$ and $t$. Therefore, the variance of $\epsilon$ for a portfolio is

$$\sigma^2(\epsilon_{Pt}) = \frac{\sigma^2(\epsilon)}{L}. \qquad (17)$$

Furthermore, $\hat{\beta}_i$ is an unbiased OLS estimator containing sampling error $e_i$. With the assumption $\sigma^2(\epsilon_i) = \sigma^2(\epsilon)$, it follows that

$$var(\hat{\beta}_i|\beta_i) = \sigma^2(e_i) = \frac{\sigma^2(\epsilon_i)}{\phi} = \frac{\sigma^2(\epsilon)}{\phi}, \qquad (18)$$

where

$$\phi = \sum_{t=1}^{T}(f_t - \bar{f})^2, \qquad (19)$$

and therefore in combination with equation (17),

$$var(\hat{\beta}_K|\beta_K) = \sigma^2(e_K) = \frac{\sigma^2(\epsilon)}{L\phi}, \qquad (20)$$

thus

$$\sigma^2(e_i) > plim_{L\to\infty}\frac{\sigma^2(\epsilon_K)}{L\phi} = 0 = \sigma^2(e_P). \qquad (21)$$

Equations (18) to (21) prove the notion that portfolio betas can be estimated more

accurately as the variance of the measurement error ($e$) diminishes for the portfolios $P$ when $N$ grows infinitely, and therefore $L$ as well, if $M$ is kept constant.

As a recap, in equation (17) Black et al. (1972) show that the variance of the error term is smaller for a portfolio. This implies that the larger the number of assets $L$ in one portfolio, the smaller the variance of the error term $\epsilon$ of the portfolio in question. Further on, in equation (18) they state that, with the assumption that the variance of the error term $\epsilon$ is constant and independently distributed, the variance of the error term $\epsilon$ divided by the variance of the explanatory variable $f$ equals the variance of $\hat{\beta}_1$. This relation is explained by Newbold, Carlson, and Thorne (2013, pp. 438-439). The properties of (17) and (18) combined (as in equation (20) and equation (21)) then lead to the conclusion that the measurement error $e$ moves asymptotically towards 0 for an increasing number of assets $L$ within one portfolio.

However, the just described logic is only half the story. In order to prove that the forming of portfolios actually improves the $\gamma$ estimate ($\hat{\gamma}$) within the second-pass regression, Black et al. (1972, p. 49) make use of equation (14) in the following rearranged form:

$$plim\ \hat{\gamma}_1 = \gamma_1 \left( \frac{\sigma_{\beta_1}^2}{\sigma_{\beta_1}^2 + \sigma_{e_1}^2} \right) = \frac{\gamma_1}{1 + \frac{\sigma_{e_1}^2}{\sigma_{\beta_1}^2}}. \tag{22}$$

Plugging in the variables and findings from equation (21), equation (22) can be stated as follows.

$$plim\ \hat{\gamma}_{1P} = \frac{\gamma_{1P}}{1 + \frac{\sigma_{e_{1P}}^2}{\sigma_{\beta_{1P}}^2}} = \frac{\gamma_{1P}}{1 + \frac{plim\frac{1}{L}\sigma^2(\epsilon)}{\phi\sigma_{\beta_{1P}}^2}}. \tag{23}$$

Thus, as shown in (21) $plim_{N\to\infty}\frac{1}{L}\sigma^2(\epsilon) = 0$ and therefore $plim_{N\to\infty}\hat{\gamma}_1 = \gamma_1$. However, Black et al. (1972, p. 49) stress the importance of the betas of the various portfolios P being dispersed as much as possible. If this were not the case then

$$plim_{N\to\infty}\sigma_{\beta_{1P}}^2 = plim_{N\to\infty}\frac{\sigma_{\beta_1}^2}{L} \tag{24}$$

would hold and cancel out $L$ in equation (23) and therefore cancel out the convergence of $\hat{\gamma}_1$ towards $\gamma_1$ for large sample sizes.

While Jagannathan, Skoulakis, and Wang (2009, p. 75) acknowledge the effectiveness of using portfolios in order to get more precise beta estimates, they also point out that one has to avoid data snooping biases as discussed by Lo and MacKinlay (1990). They point out the risk of inferences being flawed when properties of the data are used to structure

the test statistics, which is the case when the portfolios are sorted by the estimated betas.

Black et al. (1972, p. 10-11,48) were already aware of this issue and took it into account by not using estimated betas of the period under scrutiny to form portfolios. They used estimated betas from previous time periods, as these betas are assumed to be highly correlated with the betas of the time period analysed and nevertheless are independent of the measurement error $e$ (Black et al., 1972, p. 9).

However, the data snooping bias goes beyond the error-in-variables mitigation method by Black et al. (1972) and is considered to be problem for empirical research on factor models. White (2000) argues that data snooping occurs when the same data is reused for various models, eventually leading to a satisfactory result for one of the many models tested. Within this thesis, the danger of data snooping mentioned by Jagannathan et al. (2009, p. 75) should not be a concern as the model is not being adjusted within the simulation. Nonetheless, in the big picture, as White (2000) further elaborates, most of the models relying on time-series, such as the three-factor model, are probably affected by the data snooping bias. This problem is more pressing for empirical research as a whole and not for the particular examination of the intercept restriction within the three-factor model. Thus, following Black et al. (1972), the logic of the portfolio beta approach is applied within this thesis.

To sum up, portfolio betas can be measured with less error than single stocks. In a regression analysis, this effect would be offset if the portfolio betas would not vary across the formed portfolios. In other words, a high dispersion of the portfolio betas is favourable. Furthermore, it is important that the stocks are not grouped by the estimated betas of the time period under scrutiny but rather by some variable highly correlated to the beta but independent of the measurement error. For example, Black et al. (1972, p. 9) use beta estimates from past data. Moreover, well dispersed betas lead to smaller standard errors in the second-pass regression. These insights are considered in section 3, where the underlying data is presented and used for the calibration of the simulation parameters. Before that, the underlying model for the simulation approach is presented in the subsequent section.

## 2.4   The Fama French Three-Factor Model

The three-factor model introduced by Fama and French (1993) is one of the most popular multi-factor models which dominate empirical research (Cochrane, 2001, p. 437). The well-documented scientific discussion about this model is also the reason why it has been chosen for the simulation approach in this thesis. The goal of this section is to describe the characteristics of the three-factor model by contrasting it to the CAPM. As a fist step in this section, the CAPM will be briefly discussed as it "paved the way" for the three-factor model.

The introduction of the CAPM by Sharpe (1964) and Lintner (1965) ignited the empirical research of asset pricing models (Fama & French, 2004, p. 25). According to Cochrane (2001), it is "the first, most famous, and (so far) most widely used model in asset pricing" (2001, p. 152). Naturally, due to the long lasting attention given to the CAPM, it has been scrutinized by various researchers (Fama & French, 2004, pp. 30-37). Within the extensive scientific discussion about the CAPM, extensional and alternative theories and models have been introduced.

Among other reasons, the simplicity and intuitive appeal of the CAPM are key factors for its early popularity (Fama & French, 2004, p. 29). The CAPM equation is most commonly stated in equivalent return-beta language as

$$E(R_i) = R_f + [E(R_M) - R_f]\beta_{iM}, \tag{25}$$

where it describes the expected return $E(R_i)$ of asset $i$ as a composite of the risk-free rate $R_f$ and the sensitivity to the market portfolio excess return ($\beta_{iM}$) times the market excess return $[E(R_M) - R_f]$ (Cochrane, 2001, p. 152).

However, as Fama and French (2004, pp. 35-36) point out, numerous studies question the validity of the CAPM. Early tests find through

$$R_{it} - R_{ft} = \alpha_i + \beta_{iM}[R_{Mt} - R_{ft}] + \epsilon_{it}, \tag{26}$$

as a regression model for the CAPM equation (25) that the intercept is greater than the risk-free rate and that the slope of the market excess return ($\beta_{iM}$) is smaller than predicted by the model (25) (2004, p. 32).

Furthermore, Fama and French (2004) point out that in more recent tests, starting with Basu (1977) who shows that returns on stocks with a high earnings-price ratio yield higher returns than predicted by the CAPM, the focus shifts to CAPM's incapability to explain variations in these average returns (2004, pp. 34-37). Findings, such as by Rosenberg, Reid, and Lanstein (1985), who showed that the CAPM does not explain the higher than average returns of stocks with a high book-to-market equity ratio (B/M), and findings by Banz (1981), who demonstrated that average returns of stocks with small market capitalisation are higher than predicted by the CAPM, were synthesised by Fama and French (1992). This research laid the foundation for the three-factor asset pricing model, written as

$$E(R_i) - R_f = \beta_{iM}[E(R_M) - R_f] + \beta_{iS}E(SMB) + \beta_{iH}E(HML) \tag{27}$$

and commonly referred to as the Fama and French three-factor model. $SMB_t$ denotes the

"difference between returns on diversified portfolios of small and big stocks" (in terms of market capitalisation) and $HML_t$ denotes the "returns on diversified portfolios of high and low B/M stocks"(Fama & French, 2004, p. 38). The two additional betas are interpreted as exposure to the two additional factors (Cochrane, 2001, p.81).

In order to illustrate that their own factor model is able to explain much of the variation over time for portfolios with low and high B/M-ratios (value stocks) as well as portfolios with low and high market capitalisation (size), Fama and French form 25 portfolios sorted by size and value to get the required dispersion in their test data (1993, p. 10). Fama and French (1996) ran time-series regression on the excess returns of these portfolios as

$$R_{it} = \alpha_i + \beta_{iM}[R_{Mt} - R_{ft}] + \beta_S(SMB_t) + \beta_H(HML_t) + \epsilon_{it}. \tag{28}$$

If (27) perfectly captured the variation in average returns over time for each of the 25 portfolios then the pricing error $\alpha$ should result in 0 through the OLS of (28). Fama and French reject the hypothesis that $\alpha$ is equal to 0. However, they further show that, based on the same test statistics, the $\alpha$ of the three factor regression is much smaller (0.093) than the $\alpha$ of the CAPM (0.286) (1996, p. 71). In other words, the three-factor model may not be a perfect model, but it nevertheless captures a substantial part of variation of average portfolio returns and substantially more than the CAPM.

Interestingly, even though they only use time-series regression, Fama and French's (1996) inferences are about the explanatory power of the model across assets. What they show is that their model explains much of the variation of average returns for each of the 25 portfolios, which display substantially different average returns. In other words, they show that their model is capable of explaining the returns for portfolios sorted by "anomalies" that the CAPM could not explain, by achieving consistently low alphas, except for one case (Fama & French, 1996, pp. 59,71)[7]. However, it seems necessary to point out that the argument over the theoretical explanation for the empirically found relation between the size factor, the value factor and the asset return is not settled. Fama and French themselves describe the three-factor model as a "brute force construct" out of previous research which indicated anomalies not captured by the CAPM (2004, p. 39).

In short, the Fama-French three-factor model is advantageous in terms of statistical explanatory power (smaller $\alpha$) and it is widely used. However, so far, the explanatory power of the additional factors could not be explained through theoretical approaches. The theory of the three-factor model implies that $\alpha$ should be zero. Imposing this restriction on the first-pass regression as well imposing the restriction stated in formula (4) of section 2.1 on the cross-sectional regression are the foundation of this simulation study of which the methodology is outlined in the subsequent section.

---

[7]Fama and French used cross-sectional regression analysis in one of their previous study (1992).

# 3    Methodology

In order to study the effects of imposing the $H_0$ restrictions on the three-factor model within the two-pass approach, a data generating simulation has been carried out. The concept of this simulation is based on simulations done by Bänziger and Gramespacher (2015) as well as Shanken and Zhou (2007, pp. 54-57).

In a first step, the parameters for the simulation need to be calibrated. Both of the aforementioned simulation studies use the 25 book-to-market and size portfolio provided by French (2016). Once all the parameters have been calibrated, artificial returns are generated on which the first- and the second-pass regressions are run. The process of generating artificial returns and running the two regressions is repeated a number of times (e.g. 10,000 times).

In order to outline this process in detail and to identify potential problems, the rest of this section is structured as follows: first, the underlying data is presented. Then, the calibration of the simulation parameters is introduced. Finally, each of the simulation steps is described.

## 3.1    The Data

Bänziger and Gramespacher (2015) use the 5x5 value weighted book-to-market and size portfolios and the excess return of the U.S. stock market portfolio from January 1964 to October 2014. This data is continuously updated and made available by French (2016). Within this analysis, 625 monthly average returns (from January 1964 to February 2016) of the same 25 portfolios are used in order to calibrate the simulation parameters.

French (2016) constructs the data based on the procedure used by Fama and French (1992, p. 429) when they tested variants of multi-factor models. They use observed factors (value and size) from 6 months ago in order to assign the assets to the respective portfolio. Therefore, the previously outlined selection bias is avoided (see Black et al., 1972, p. 9). By sorting the assets into portfolios based on these historically observed factors, they create 25 portfolios with dispersed average returns, as they account for the "anomalies" found by the various researchers mentioned in section 2.4. The motivation of Fama and French for this portfolio sorting is that they want to have test statistics displaying the "anomalies" not explained by the CAPM, as they want to show that their three-factor model captures more of the variations of average returns across assets (1996, p. 68).
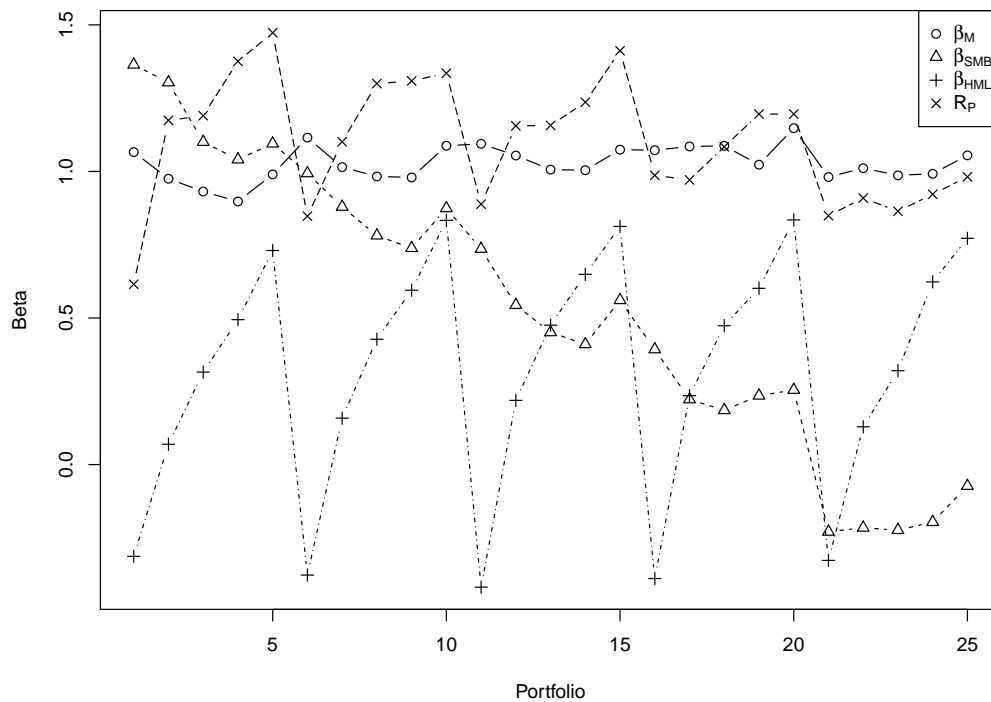
In short, the fact that Bänziger and Gramespacher (2015) use the same data set as well as the fact that the data set is well known (Shanken & Zhou, 2007, p. 54) are the reasons why they are used within this simulation. In addition to the data used by Bänziger and Gramespacher (2015), the two additional monthly observed factors are retrieved from the

database provided by French (2016).

## 3.2 The Parameters

With the data introduced in section 3.1, the parameters for the simulation are now calibrated by using the whole observation period. This includes the statement of the true values of the three factors, the estimation of the betas for the three risk-factors, the estimation of the covariance of the error terms resulting from the beta estimations, as well as the calibration of the covariance of the three factors. Each of these parameters are now described in detail within this section.

The true value for the factor risk premium are $\gamma_1 = 0.4824$ for the market factor, $\gamma_2 = 0.2349$ for the size factor (SMB) and $\gamma_3 = 0.3404$ for the value factor (HML). These values are simply the mean of the 625 monthly three-factors observed by French (2016).



**Figure 1:** Three factor betas dispersion across 25 portfolios.

The estimated betas are illustrated in figure 1. It shows the dispersion of the estimated betas for the three factors over the whole observation period. These betas are estimated with an OLS time-series regression through the origin with the excess returns of the 25 portfolios as the regressand and the three factors as the regressor. Leaving out the intercept for estimating the population betas is done in accordance with Shanken and Zhou (2007, p. 55). Figure 1 indicates that the market-beta $\beta_M$ does not capture much of the variation of

the average portfolio returns ($R_P$) across the 25 portfolios, as it does not vary substantially even though the average returns vary across the portfolios. In other words, the market-beta explains a part of the portfolio return (unless $\gamma_1 = 0$, which is not the case), but only little of the variation of average returns across the portfolios. The two additional betas of the three-factor model on the other hand, notably HML, seem to capture much of the variation of the average market return across the 25 portfolios.

However, figure 1 also indicates a potential obstacle for the accuracy of the two-pass method. In the analysis by Bänziger and Gramespacher (2015), where $\beta_M$ is the sole factor for capturing the relation between the excess return of portfolio $i$ and the excess return of the market, the dispersion of this factor is substantial (from 0.8 to 1.4) (2015, p. 79). Within the present analysis, most of the variation in $\beta_M$ is gone - it seems as though the two additional factor capture the variation across the average portfolio returns (and more) previously captured by $\beta_M$. Thus, now that $\beta_M$ hardly varies across the 25 portfolios, it seems reasonable to expect stronger effects due to the error-in-variables problem outlined in section 2.3, notably the problem outlined in equation (23) and equation (24) about the much needed dispersion of the portfolio betas within the second-pass regression.

Through the estimation of the betas, error terms for each of these regressions are obtained. These are used to estimate the covariance matrix $E$ of these error terms. Including this parameter in the simulation allows to account for problems caused by correlation of the residuals in the cross-sectional regression (Cuthbertson & Nitzsche, 2004, pp. 201-202). Fama and MacBeth (1973) were among the first to address this problem with their "rolling-regression" approach (Fama & French, 2004, p. 31). Fama and French (1992, p. 39) themselves used a data generating simulation approach also capturing the effects of residual correlation of cross-sectional regressions. In this sense, Bänziger and Gramespacher (2015) use a similar approach to Fama and French (1992), which will elaborated in detail in section 3.3.

Finally, the covariance of the three factors is derived. This parameter is required to generate the artificial returns within the simulation. As is described in section 3.3, the generation of the returns is done by assuming that the three-factor model holds. Therefore, when generating returns by applying the factor model, data sets of these three factors are required to resemble the true data.

To sum up, the parameters for the data generating simulation are the true means of the risk factors, the estimated betas of these factors for every portfolio, the estimated cross-sectional covariance of the error-terms and the covariance of the three-factors among themselves. All these parameters are derived out of the whole data set (625 monthly observations) and form the foundation of the simulation outlined in the subsequent section.

## 3.3   The Simulation

The simulation is programmed in the statistical programming language $R$ and can essentially be divided into four parts. First, returns of the sample size $T$ are generated. Second, the returns are used as the regressand within the first-pass regression. Third, the coefficients resulting from the first regression are used within the second-pass regression as the explanatory variables and the returns as the regressand. These first three steps are then repeated 9,999 times. Finally, the statistics are drawn from the 10,000 results. This simulation can then be repeated for various sample sizes $T$. After a brief explanation on the sizes chosen for $T$ as well as the underlying assumptions concerning the distribution of the artificially generated returns, the rest of this section is structured along the four aforementioned steps. A detailed description of the $R$-code is available in the appendix.

As mentioned, the $\beta_M$ properties of the calibrated data are somewhat problematic. Therefore, in addition to the sample sizes $T_1 = 60$ and $T_2 = 120$ chosen by Bänziger and Gramespacher (2015), data sets for a sample size of $T_3 = 360$ are simulated as well, since the problem of error-in-variables should diminish with larger sample sizes (Shanken & Zhou, 2007, p. 55). However, it is questionable whether the implied assumption that betas are stationary over 360 months is appropriate (Bänziger & Gramespacher, 2015, p. 78). Nonetheless, this additional loop of simulations should provide statistics to put the expected error-in-variables problem of $\beta_M$ into perspective. For each $T$ specification $(T_1, T_2, T_3)$ 10,000 data sets are drawn.

Furthermore, for the generation of artificial returns, it is assumed that the returns are normally distributed since Shanken and Zhou (2007, p. 55) as well as Bänziger and Gramespacher (2015, p. 79) generate normally distributed returns within their simulation. This assumption is regarded as a standard assumption (Shanken & Zhou, 2007, p. 55). In addition, as done by Bänziger and Gramespacher (2015), it is assumed that the observed returns are not serially-correlated and that the error variance is not related to the explanatory variables (assumption of homoscedasticity). Cochrane (2001, p. 285) provides evidence that simulating returns which do not account for the heteroscedasticity and auto-correlation of the underyling data is less of a problem when analyzing a factor model through a two-pass approach[8]. To put it differently, since the simulation by Bänziger and Gramespacher (2015) did not take into account these properties of the underlying data and since there is evidence that the consequences are limited, the same assumptions are made. Conversely, the cross-sectional correlation of the error terms will be considered when simulating the artificial returns for the regression analysis.

In order to generate $T$ returns, $T$ sets of the three factors are required. These can be generated by making use of the *mvrnorm* function available in the $R$ library *MASS*. The

---

[8]Cochrane (2001, pp. 235-236) refers to the two-pass method as simply the cross-sectional regression.

inputs are the true values of the three factors ($\gamma_1, \gamma_2, \gamma_3$ as mentioned in section 3.2) as well as the covariance among themselves over all 625 observations. The resulting $Tx3$ matrix $F$ contains $T$ normally distributed sets of the three factors. $F$ is then used to simulate $T$ returns for the 25 portfolios.

Moreover, since the returns to be generated are supposed to reflect the properties of the true data, the previously estimated covariance of the error term is taken into account. Since the error terms are assumed to be independent over time, the artificial error terms can be generated through the use of the *mvrnorm* as well. The resulting $TxN$ matrix $S$ is created with the estimated covariance of the residuals $E$ and a mean of 0 [9]. However, by imposing the assumption of homoscedasticity and non-serial correlation, the test statistics do not reflect all properties of the observed data. The data generating equation takes the form

$$FB' + S = R, \tag{29}$$

where $B$ is a $25x3$ matrix and contains the estimated betas when estimated for the whole observation period (as depicted in figure 1). These generated returns $R$ are now being used within the two-pass approach. A detailed description of the data generating process is provided in the appendix. Once the returns have been simulated, the two regressions can be run on these returns.

With the three-factor model as the underlying model, the first-pass regression takes the form

$$R_{i,t} = \alpha_i + \beta_{i1}f_{1t} + \beta_{i2}f_{2t} + \beta_{i3}f_{3t} + \epsilon_{it}, \qquad i = 1, ..., N, \qquad t = 1, ..., T, \tag{30}$$

where $f_{1t}$ = return of the market portfolio in excess to the risk-free rate, $f_{2t}$ = the size factor (SMB) and $f_{3t}$ = the value factor (HML) of period $t$. All other variables are defined as in (1). Consequently, the second-pass regression takes the form

$$\bar{R}_i^e = \gamma_0 + \gamma_1\hat{\beta}_{i1} + \gamma_2\hat{\beta}_{i2} + \gamma_3\hat{\beta}_{i3} + u_i, \tag{31}$$

where all variables are defined as in section 2.1.

As two methods are applied for the first-pass regression, two sets of $T$ beta estimates (loading factors) are generated in the first-pass regression. Thereafter, the resulting two sets of loading factors are run in the second-pass (cross-sectional) regression, again with and without the intercept restriction. Thus, four sets of risk premium estimations are generated for every loop. After storing the derived results in a matrix, the whole process is

---

[9]As stated in equation (1) $E(\epsilon_{i,t}) = 0$.

looped 9,999 times. Once this process is finished, the averages of the 10,000 results are compared with the true values in order to examine whether the various estimation methods perform well. A detailed description of this process is provided in the appendix.

As outlined in section 3.2, the dispersion of the two additional betas in figure 1 might lead to small biases. However, it might not be desirable to have small error-in-variables biases when the goal is to study possible mitigation effects for them. Or to put it differently, in order to study the effect on the bias, there has to be a bias. Therefore, in addition to the standard simulation, one run for the $T = 60$ sample size is carried out with the estimated covariance across the residuals multiplied by the factor 100. The magnitude of this factor is chosen arbitrarily. This should lead to higher measurement errors, since the residual variance is the cause for the measurement error of the portfolio betas (Cochrane, 2001, p. 434).

To sum up, within this simulation returns are generated that reflect the true properties of the observation period in terms of cross-sectional correlation across the error-terms of the beta estimates. In addition, the covariance of the three factors among themselves is also taken into account within the data generating process. The data is generated with the assumption that the returns and error terms are normally distributed. The error terms are assumed to be independent over time but not across assets. The generated returns and the underlying calibrated parameters are the foundation to run the first- and second-pass regression multiple times. This set-up enables the analysis of the effects of imposing the restriction that $\alpha = 0$ within the first-pass regression and the restriction $\gamma_0 = 0$ within the second-pass regression. In the following section, the results of this simulation approach are presented.

# 4   Simulation Results

This section is structured in the following way: in the first part, the findings are summarized. For this purpose, all the estimation results are presented in a tabular structure. Then, in order to verify whether the simulation results are reasonable, the results of the case where the first- and second-pass regressions are run with intercepts is analyzed. Comprehensive results for this method are made available by Shanken and Zhou (2007). By comparing the results to their study, it is discussed whether the obtained results in this thesis are reasonable for the method where no intercepts are omitted . Further on, the extended simulation results (regression through the origin in at least one of the regressions) are displayed and contrasted to both Bänziger and Gramespacher (2015) and Shanken and Zhou (2007). In a second part, the findings are discussed.

## 4.1   Findings

The results of the simulation are displayed in table 1. It shows the means of all 10,000 estimates for the intercept ($\gamma_0$), the market premium ($\gamma_1$), the size premium ($\gamma_2$), the value premium ($\gamma_3$), as well as all the corresponding standard deviations of these 10,000 estimated coefficients. The simulation was carried out for three different sample sizes ($T$=60, $T$=120, $T$= 360). The column "Method" states the methods for the first-pass regression (before the comma in the abbreviation) and the second-pass regression (after the comma in the abbreviation). "w/" means with intercept and "w/o" means without intercept.

The method "w/, w/" in the first row of table 1 corresponds to the method applied by Shanken and Zhou (2007, p. 64). The estimation results within this thesis are similar to their results in terms of direction of the bias and asymptotic behaviour for growing sample sizes $T$. The bias of the estimated premiums is the strongest for the smallest sample size ($T$=60), of which the market-premium ($\gamma_1$) contains the most pronounced downward bias. The average of the 10,000 $\gamma_1$ estimates is 42% lower than the true value. The other two premiums are underestimated by 3.04% ($\gamma_2$) and 4.29% ($\gamma_3$). These discrepancies decrease with an increasing sample size $T$. However, even when the sample period is 15 years ($T$ = 120), the estimated intercept is still substantially greater than zero.

Unfortunately, Shanken and Zhou (2007) do not fully state their numeric results for their cross-sectional estimations in the case of three factors. However, they do display the results for the one-factor case (CAPM) and find that the intercept estimation as well as the market-premium are qualitatively equal to the findings of the one-factor case. They observe that the biases are larger in the case of three factors and the standard deviation of the estimates are smaller than in the case of the CAPM (Shanken & Zhou, 2007, p. 63). The results presented in the first row of table 1 feature the just mentioned properties when compared to the one-factor results of Shanken and Zhou (2007, p. 56) as well as to the results of Bänziger and Gramespacher (2015, p. 80).

**Table 1:** Results of Four Estimation Methods

| Method | Estimate | Length of Time Series | | | | | |
|---|---|---|---|---|---|---|---|
| | | T = 60 months | | T = 120 months | | T = 360 months | |
| | | Mean | SD | Mean | SD | Mean | SD |
| w/ , w/ | $\gamma_0$ | 0.21713 | 0.8021 | 0.13112 | 0.62802 | 0.04901 | 0.3865 |
| | $\gamma_1$ | 0.27836 | 0.8699 | 0.37019 | 0.70660 | 0.43562 | 0.4443 |
| | $\gamma_2$ | 0.22775 | 0.4086 | 0.23590 | 0.28711 | 0.23301 | 0.1663 |
| | $\gamma_3$ | 0.32580 | 0.3708 | 0.33410 | 0.26272 | 0.33832 | 0.1538 |
| w/o , w/ | $\gamma_0$ | -0.00425 | 0.7180 | -0.00048 | 0.5867 | 0.00018 | 0.3683 |
| | $\gamma_1$ | 0.48377 | 0.9177 | 0.49282 | 0.7148 | 0.47965 | 0.4369 |
| | $\gamma_2$ | 0.23972 | 0.4130 | 0.24081 | 0.2883 | 0.23467 | 0.1664 |
| | $\gamma_3$ | 0.34095 | 0.3780 | 0.34277 | 0.2650 | 0.34119 | 0.1540 |
| w/ , w/o | $\gamma_1$ | 0.48368 | 0.5781 | 0.49438 | 0.4134 | 0.48201 | 0.2379 |
| | $\gamma_2$ | 0.24054 | 0.4103 | 0.23954 | 0.2884 | 0.23434 | 0.1669 |
| | $\gamma_3$ | 0.33530 | 0.3707 | 0.33708 | 0.2624 | 0.33945 | 0.1537 |
| w/o , w/o | $\gamma_1$ | 0.47970 | 0.5791 | 0.49240 | 0.4135 | 0.48140 | 0.2377 |
| | $\gamma_2$ | 0.24187 | 0.4144 | 0.24086 | 0.2894 | 0.23473 | 0.1668 |
| | $\gamma_3$ | 0.34024 | 0.3779 | 0.34271 | 0.2649 | 0.34122 | 0.1540 |

Note: Monthly values in percent - abbreviation "w/" stands for regression with intercept. Abbreviation "w/o" stands for regression without intercept. The abbreviation before the comma states the method for the first-pass regression, the latter the method for the second-pass regression. For each sample size length, $T = 10,000$ data sets were generated. SD = standard deviation. The true values: $\gamma_0 = 0$, $\gamma_1 = 0.4824$, $\gamma_2 = 0.2349$ and $\gamma_3 = 0.3404$.

Shanken and Zhou (2007, pp. 64-65) display their results for the value premium and state that the size premium estimation results were qualitatively the same. The results are similair to the results of this thesis for the value premium ($\gamma_3$) estimated with the "w/, w/" method. The relative underestimation of the value premium in the "w/, w/" method with $T = 60$ is 0.04% greater than the one of Shanken and Zhou (2007, p. 64). However, the absolute comparison of these values is probably flawed, as Shanken and Zhou (2007, p. 55) use substantially different calibration values for their simulation and observe a slightly shorter time period. Furthermore, they do not give the numerical results of the size premium, but state that they are qualitatively similar to the value premium (Shanken & Zhou, 2007, p. 63). This is also the case for the results shown in table 1.

To sum up, the estimations are biased for the "w/, w/" method. The bias is the strongest for $\gamma_1$. The bias diminishes with larger sample sizes $T$. The results are in line with the OLS method by Shanken and Zhou (2007, pp. 54-64), which is considered a strong indication

that the simulation within this thesis is well-defined. With this insight, the results of the extensional methods are now be presented.

Figures 2 through 7 illustrate the rest of the results stated in table 1. They are displayed in pairs. Each pair presents the results for one type of risk premium. Figures 2 and 3 display the estimations of $\gamma_1$, figures 4 and 5 display the estimations of $\gamma_2$ and figures 6 and 7 display the estimations of $\gamma_3$. The figures on the left of each pair display the results for the sample size $T = 60$. Consequently, the three figures on the right display the results for the sample size $T = 120$. Each figure depicts the estimations for the four methods. The results form the $T = 360$ simulation run are left out, as its main purpose was to verify whether the simulation would behave, as expected for large sample sizes, which is the case. Finally, the results for the simulation run with a distorted covariance estimation of the residuals (multiplied by the factor 100) are graphically displayed in figures 8 and 9. The numeric results for the latter are not shown. The rest of the findings presentation follows the order of these 8 depictions.
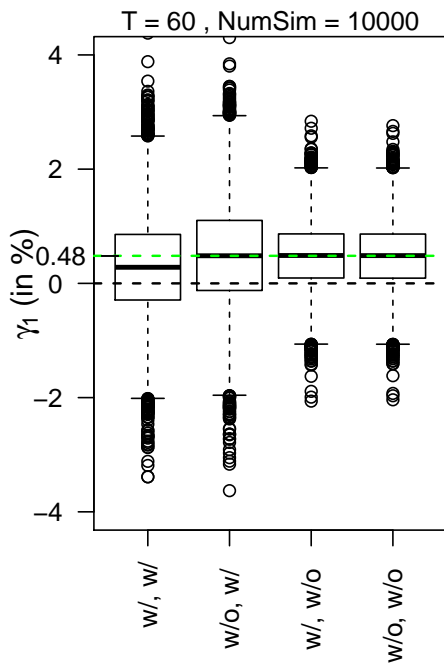


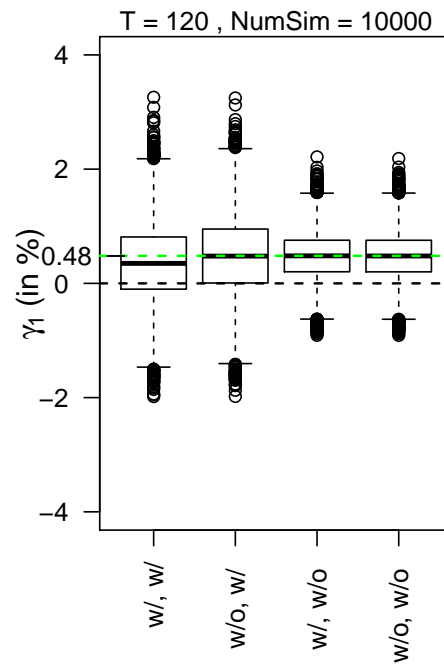**Figure 2:** $\gamma_1$ estimate, sample size $T = 60$.     **Figure 3:** $\gamma_1$ estimate, sample size $T = 120$.

Figures 2 and 3 illustrate three things prominently: first, they highlight the strong downwards bias of $\gamma_1$. Second, similar to Bänziger and Gramespacher (2015, p. 80), forcing at least one of the regression through the origin decreases the bias of the average simulation estimate substantially. And finally, estimating the factor-risk premium without an intercept in the second-pass regression decreases the standard deviation of the estimate
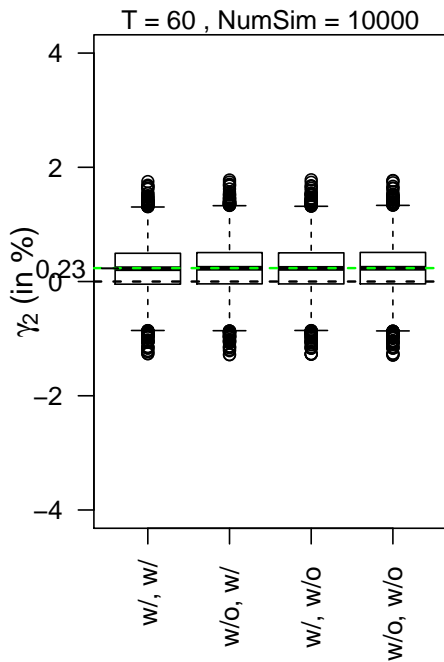
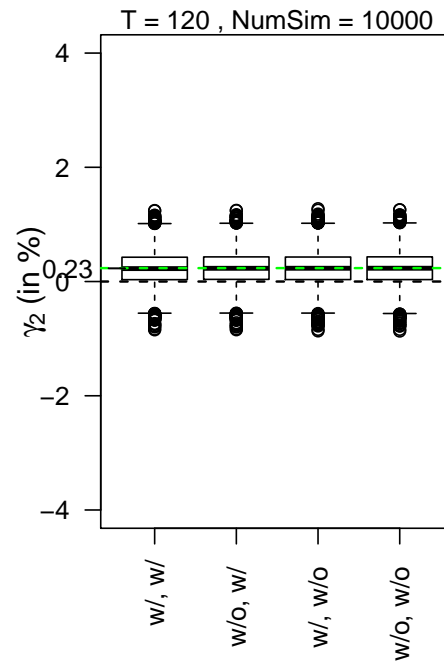**Figure 4:** $\gamma_2$ estimate, sample size $T = 60$.

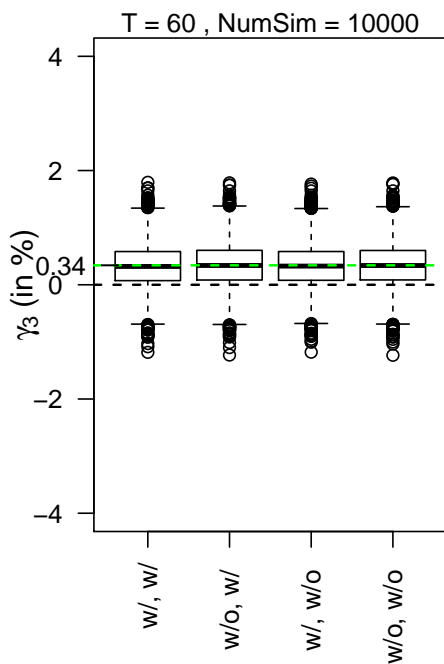**Figure 5:** $\gamma_2$ estimate, sample size $T = 120$.



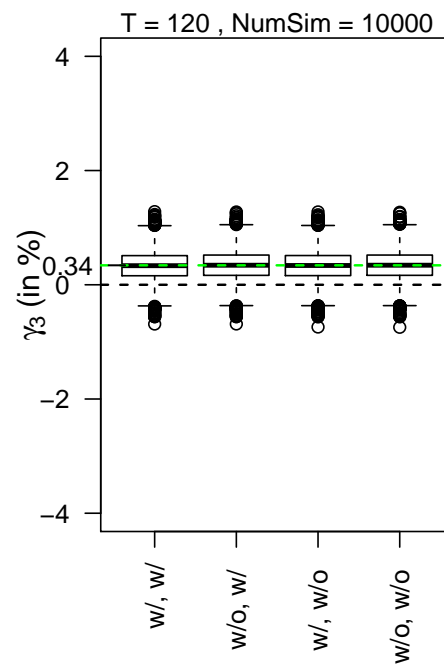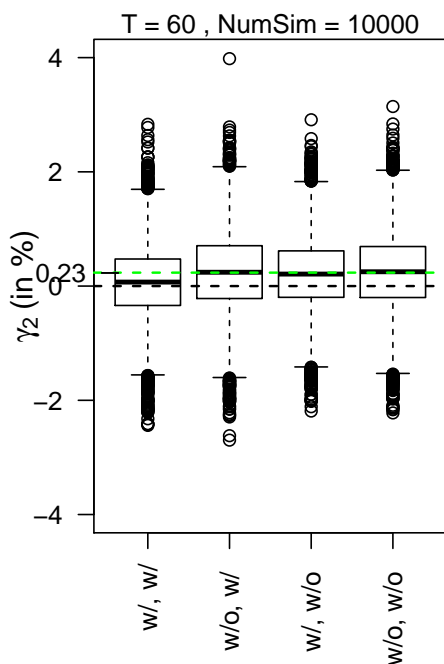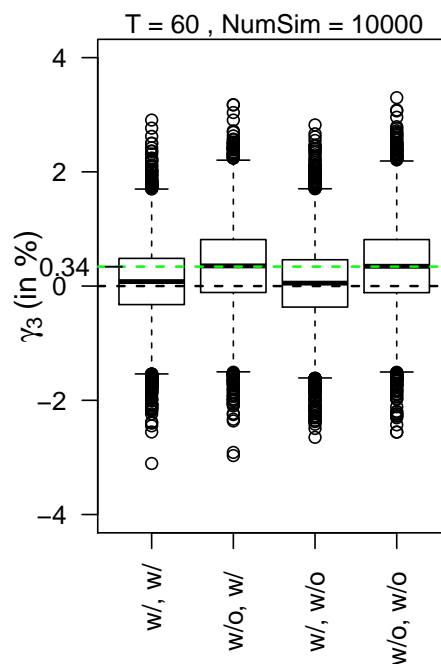**Figure 6:** $\gamma_3$ estimate, sample size $T = 60$.

**Figure 7:** $\gamma_3$ estimate, sample size $T = 120$.

**Figure 8:** $\gamma_2$ estimate, sample size $T = 60$ with distorted residual covariance.

**Figure 9:** $\gamma_3$ estimate, sample size $T = 60$ with distorted residual covariance.

substantially. In other words, the mitigating effects shown by Bänziger and Gramespacher (2015) also occur for the market risk premium, even when it is not the only explanatory factor within the second-pass regression.

Figures 4 to 7 highlight similar patterns for $\gamma_2$ and $\gamma_3$. The four figures show that the standard deviation of the estimates diminishes with larger sample sizes ($T$). However, these premium estimations seem to be indifferent to forcing the first- or second-pass regression through the origin. No clear reduction of a downward bias can be observed nor does leaving out the intercept in the cross-sectional regression lead to a smaller standard deviation for the estimates. However, the numeric display of the results in table 1 show small reductions in the bias when the constant in the first- or second-pass regression is omitted. This could be an indication that if the bias of the two additional premium estimates would be larger the mitigating effect on the bias would be clear. In contrast, the numeric results in table 1 do not yield evidence for a mitigating effect on the standard errors of the two additionally estimated premiums.

Figures 8 and 9 depict the case of exaggerated cross-sectionally correlated error terms. Strong limitations for this extra simulation are given in the discussion and the conclusion. Nonetheless, it indicates that when a clear bias exists, leaving out the intercept in the cross-sectional regression leads to substantially more accurate estimations. However, the estimates of $\gamma_3$ does not improve within the method "w/, w/o" whereas this is the case for

the estimates of $\gamma_2$. Furthermore, no indication is visible that leaving out the intercept in the second-pass regression leads to smaller standard errors of the estimates.

In short, the results indicate that omitting the intercept in at least one of the two regressions decreases the bias of the $\gamma_1$ estimates. Furthermore, leaving out the intercept in the second-pass regression leads to a smaller standard deviation for the $\gamma_1$ estimates. The mitigating effect of omitting the constants in at least one of the two regressions is less clear for the size and value premium estimation, since the estimation bias of these two premiums is already small without omitting the intercept in either of the regressions. Moreover, the standard errors of the size and value premium are not reduced when omitting the intercept in at least one of the two regressions. Lastly, in accordance with the studies of Bänziger and Gramespacher (2015) as well as Shanken and Zhou (2007), the bias in the estimates and the standard errors of the estimates diminish for larger sample sizes.

## 4.2   Discussion of the Simulation Results

This simulation is an extension to the simulation by Bänziger and Gramespacher (2015) in terms of factor and it is an extension to the simulation by Shanken and Zhou (2007) in terms of estimating the first- and second-pass regression without constants. In order to contrast the findings of this thesis with the two aforementioned studies, this section is structured in the following way: first, the "w/, w/" method is discussed. Then, the findings are discussed in the order of figure 2 to figure 9.

The comparison between the "w/, w/" results and the results of Shanken and Zhou (2007) indicate that the simulation is set-up properly. The results vary in the same way for different sample sizes $T$ and the relative magnitude of the bias are alike. For example, the downward bias of $\gamma_1$ is by far the most serious bias and all the biases diminishes for larger sizes of $T$. Furthermore, the comparison of the results to the results of Bänziger and Gramespacher (2015) shows that the bias of $\gamma_1$ is substantially larger for the three-factor model than for the CAPM, which is in line with the findings of Shanken and Zhou (2007). However, the similar results for the same method are no proof that the simulation is well specified. Nonetheless, it is considered a strong indication for this notion.

The estimation results of $\gamma_1$ (depicted in figures 2 and 3) are qualitatively similar to the findings of Bänziger and Gramespacher (2015). The market premium estimation bias decreases for larger sample sizes $T$ and leaving out the intercept in one of the two regressions leads to substantially more accurate estimates for $\gamma_1$. Furthermore, if the cross-sectional regression is estimated without a constant, then the standard deviation of the estimates of $\gamma_1$ decrease considerably.

An explanation for the stronger downwards bias of $\gamma_1$ within the three-factor model might be the insufficient dispersion of the market beta shown in figure 1. As mentioned

earlier, Black et al. (1972, p. 49) stressed the importance of the betas of the various portfolios $P$ being dispersed as much as possible, as otherwise forming portfolios and generating more accurate beta measures for the individual portfolios are in vain. Within this simulation, there is still some variation across the market betas of the portfolios; however, the reduced variation might have caused the bigger downward bias for the $\gamma_1$ estimate. However, with the logic of the little dispersion, the question arises why the standard deviation of the market premium estimates are not larger than in the results of Bänziger and Gramespacher (2015), since by the logic of Black et al. (1972), outlined in section 2.3, one might expect the standard errors to be larger when the dispersion of the beta in question is lower. This thesis does not provide an answer on this. A possible explanation could be that within this simulation the interdependence of the three factor might have mitigated this problem.

Yet again, these elaborations are solely an attempt to explain the results rather than a proof. In other words, as shown in section 2.3, the lower dispersion most likely has an impact on the directional bias, but other effects are not ruled out by this. Overall, the estimation results for $\gamma_1$ behave similarly to the results of Bänziger and Gramespacher.

The case for the estimation of the two additional risk premiums (depicted in figures 4 to 7) is different. The results do not show a mitigating effect on the standard error of the estimates by running the cross-sectional regression without an intercept. Furthermore, the error-in-variables bias is already hardly observable for the "w/, w/"; therefore, the observation whether omitting the intercept in at least one of the regressions leads to less biased estimates of $\gamma_2$ or $\gamma_3$ is less clear. It seems as though the trade-off described by Cochrane (2001, p. 26) between efficiency and robustness does not apply for these two additional factors. In other words, the reduction in robustness by imposing the restriction that the constant should be zero does not yield more efficient estimates.

Figures 8 and 9 illustrate that leaving out the intercept in the second-pass regression decreases the bias of the estimates in a similar fashion to the case of the estimations for $\gamma_1$ when the bias is stronger. However, the case of $\gamma_3$ with the "w/, w/o" estimation method does not match the other observations. A possible explanation approach could be the fact that only one simulation parameter has been meddled with. Therefore, many interdependencies of the various parameters and estimations are distorted with impacts not detected in this simulation. At this point, no explanation has been found for this different estimation result between the two risk premiums. In general, this additional simulation bases on a extreme manipulation of one factor, therefore these additional results can at most be considered as an indication.

# 5  Conclusion

Within this simulation, further evidence is provided on the effects of running the first-
and/or the second-pass regression of the two-pass method without a constant. The simu-
lations carried out by Bänziger and Gramespacher (2015) are extended with two additional
risk factors. Instead of the CAPM, the Fama and French three-factor model has been esti-
mated. Furthermore, the simulation can also be considered an extension to the simulations
made by Shanken and Zhou (2007). In fact, one of the methods simulated is very close to
one of their simulations. The comparison of the results of the method used in this thesis
with the results of Shanken and Zhou (2007) indicates that the simulation is well specified,
at least for the method in question.

Similar to Bänziger and Gramespacher (2015), the estimation results suggest that leav-
ing out the intercept in at least one of the regressions decreases the bias of the market-risk
premium estimates, as well as their standard deviation. However, the results for the size
and value factor premiums are less clear. The biases of these two factors seem to be
substantially smaller than the bias of the market-risk premium. Since the average of the
estimates is already close to the true value (small bias) for the two additional factors, the
additional simulation methods did not yield clear evidence for a smaller bias in the estima-
tion of the size and value premium when leaving out intercepts. Moreover, the standard
deviation of the size and value premium estimates does not decrease when omitting the
intercept in either of the two regressions.

In order to retrieve more pronounced results for the mitigating effect on the bias, future
research could achieve stronger biases in their test statistics by using different underlying
data. As Cochrane (2001, p. 443) points out, there is no particular reason for choosing the
portfolios used within this thesis despite its popularity. As outlined in section 2.4, a factor
model should be able to explain any kind of dispersion of average returns across different
portfolios. Different data might lead to less well dispersed betas for the three factors and
therefore might lead to stronger biases in the estimates of the size and value risk premium.
This would make meddling with the simulation parameters in order to get stronger biases
obsolete.

As for the standard error reduction, one might consider trying to explain why the stan-
dard deviation of the $\gamma_2$ and $\gamma_3$ estimates do not decrease when the cross-sectional regres-
sion is run without a constant. The answer might lie in the fact that the two additional
factors were, unlike the market premium, found empirically (as outlined in section 2.4).
Therefore, by intuition, it could be questionable whether factors which were found when
empirical researchers "hunted" for factors that fit a model with the $H_0$ assumption of $\alpha = 0$
are even sensitive to imposing this restriction.

All in all, this thesis provides evidence, that leaving out the intercepts within the two-pass approach is hardly beneficial in terms of reducing the bias in the estimates for the size factor premium or the value factor premium, since the biases are already relatively small without omitting the intercepts. The theoretical attempt to explain the smaller biases is that the data used within this simulation exhibits well dispersed betas for the two additional factors, which is beneficial to the accuracy of their measurement within the simulation. Furthermore, omitting the intercepts does not lead to smaller standard errors of the estimates for the size and value premium. Even if leaving out the intercept is not necessary in order to get better measurements for the two additional factors, in the big picture, omitting the intercept when estimating the three-factor model through a two-pass approach could be a viable option, as the bias of the market factor premium estimate decreases substantially for small sample size estimates. However, the trade-off mentioned by Bänziger and Gramespacher (2015) between efficiency and robustness still applies. In fact, the "deal" became worse in the case of the three-factor model, since the estimates of the two additional premiums became less robust without any clear efficiency gain. In the end, similar to Bänziger and Gramespacher (2015), the standard errors of the estimates are substantial in all of the methods applied. Therefore, drawing statistical conclusions from the two-pass method continues to be a delicate task.

# References

Banz, R. W. (1981). The relationship between return and market value of common stocks. *Journal of Financial Economics*, *9:1*, 3–18.

Bänziger, A., & Gramespacher, T. (2015). Estimating beta pricing models with or without an intercept: Results from simulations. *International Research Journal of Finance and Economics*, *140:1*, 76-82.

Basu, S. (1977). Investment performance of common stocks in relation to their price-earnings ratios: A test of the efficient market hypothesis. *The journal of Finance*, *32:3*, 663–682.

Black, F. (1993). Beta and return. *The Journal of Portfolio Management*, *20:1*, 8–18.

Black, F., Jensen, M. C., & Scholes, M. S. (1972). The capital asset pricing model: Some empirical tests. In M. C. Jensen (Ed.), *Studies in the theory of capital markets* (p. 79-121). New York: Praeger.

Blume, M. E. (1970). Portfolio theory: A step toward its practical application. *The Journal of Business*, *43:2*, 152–173.

Brooks, C. (2014). *Introductory econometrics for finance*. Cambridge: Cambridge university press.

Cochrane, J. H. (2001). *Asset pricing*. Princeton and Oxford: Princeton University Press.

Cuthbertson, K., & Nitzsche, D. (2004). *Quantitative financial economics: stocks, bonds and foreign exchange*. Chichester: John Wiley & Sons Ltd.

Fama, E. F., & French, K. R. (1992). The cross-section of expected stock returns. *the Journal of Finance*, *47:2*, 427–465.

Fama, E. F., & French, K. R. (1993). Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics*, *33:1*, 3–56.

Fama, E. F., & French, K. R. (1996). Multifactor explanations of asset pricing anomalies. *The Journal of Finance*, *51:1*, 55–84.

Fama, E. F., & French, K. R. (2004). The capital asset pricing model: Theory and evidence. *Journal of Economic Perspectives*, *18:3*, 25–46.

Fama, E. F., & MacBeth, J. D. (1973). Risk, return, and equilibrium: Empirical tests. *The Journal of Political Economy*, *81:3*, 607–636.

French, K. R. (2016). *25 portfolios formed on size and book-to-market (5 x 5)*. Retrieved March 3, 2016, from `http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html`

Jagannathan, R., Skoulakis, G., & Wang, Z. (2009). The analysis of the cross-section of security returns. In Y. Ait-Sahalia & L. P. Hansen (Eds.), *Handbook of financial econometrics* (p. 73-134). Amsterdam: Elsevier Science.

Kennedy, P. (2008). *A guide to econometrics*. Malden, MA: Blackwell Publishing.

Lintner, J. (1965). The valuation of risk assets and the selection of risky investments in stock portfolios and capital budgets. *The Review of Economics and Statistics*, 13–37.

Lo, A. W., & MacKinlay, A. C. (1990). Data-snooping biases in tests of financial asset pricing models. *Review of Financial Studies*, *3:3*, 431–467.

Merton, R. C. (1973). An intertemporal capital asset pricing model. *Econometrica: Journal of the Econometric Society*, *41:5*, 867–887.

Newbold, P., Carlson, W., & Thorne, B. M. (2013). *Statistics for business and economics*. Harlow: Pearson Higher Ed.

Rosenberg, B., Reid, K., & Lanstein, R. (1985). Persuasive evidence of market inefficiency. *The Journal of Portfolio Management*, *11:3*, 9–16.

Shanken, J. (1992). On the estimation of beta-pricing models. *Review of Financial studies*, *5:1*, 1–33.

Shanken, J., & Zhou, G. (2007). Estimating and testing beta pricing models: Alternative methods and their performance in simulations. *Journal of Financial Economics*, *84:1*, 40–86.

Sharpe, W. F. (1964). Capital asset prices: A theory of market equilibrium under conditions of risk. *The Journal of Finance*, *19:3*, 425–442.

Van der Wijst, N. (2013). *Finance: a quantitative introduction*. Cambridge: Cambridge University Press.

White, H. (2000). A reality check for data snooping. *Econometrica*, *68:5*, 1097–1126.

Wooldridge, J. (2013). *Introductory econometrics: A modern approach*. Boston, MA: Cengage Learning.

# Appendix: R-Code Simulation

In this section, the $R$-code of the simulation is described in detail. It is the complete simulation code. First, the input format of the data as well as the required inputs are briefly described. Next, the initialization of the simulation properties as well as the calibration of the parameters is outlined. In a third step, the data generating process is displayed in detail, followed by the description of the two regression steps. Finally, for the sake of completeness, the code utilized for visualizing the results is added. However, since the focus lies on the simulation part, the last part is only accompanied by in-code comments.

As the underlying data, the 5x5 portfolio return data[10] as well as the corresponding factor data[11] for the three-factor model provided by French (2016) is used. In order to get the code working, one should store a CSV table of the 5x5 portfolio data which solely contains the observation dates (e.g. "012016" for January 2016) in the first column and the return data of the $N(=25)$ portfolios in columns 2 to 26. It is important that the headers be included and equal to the headers chosen by French (2016). The code below indicates what simulation inputs have to be chosen. $T$ denotes the sample size to be drawn $NumSim$ times. For example, within this thesis 10,000 data sets have been generated for three sizes of $T$. It is important that

$$T > (N + K)$$

holds, else the matrix inversions such as on code-line 116 will return an error.

**Listing 1:** Data and Initializing

```
1  #***********************************#
2  #Data:
3  #***********************************#
4
5  environmentName("Global Environment")
6  # Dataset: 01 - 1964 - 02 2016
7  # 25 Portfolio Data retrieved from http://mba.tuck.dartmouth.edu/pages/
       faculty/ken.french/ftp/25_Portfolios_5x5_CSV.zip
8  # The downloaded file contains more data than necessary, all data
       should be deleted apart from the "Average Equal Weighted Returns --
       Monthly" for each of the portfolios and the timeline in the first
       row and the headers for the 25 portfolios are kept.
9  Data5x5 <- read.csv("Data/25_Portfolios_5x5.csv",TRUE)
```

---

[10]http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/ftp/25_Portfolios_5x5_CSV.zip

[11]http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/ftp/F-F_Research_Data_Factors_CSV.zip

```r
10
11  #Data retrieved from http://mba.tuck.dartmouth.edu/pages/faculty/ken.
        french/ftp/F-F_Research_Data_Factors_CSV.zip
12  #The same is repeated for the 3 Factor, all data should be deleted
        apart from the monthly factors, the time line in the first row and
        the headers of the three-factor columns.
13  #If the header of the each column are stored with the original column
        headers then:
14  #Data3x3$Mkt.RF refers to the 625 observations of the market factor,
15  #Data3x3$SMB refers to the 625 observations of the size factor and
16  #Data3x3$HML refers to the 625 observations of the value factor.
17  Data3x3 <- read.csv("Data/F-F_Research_Data_Factors.csv",TRUE)
18
19
20  #*************************************#
21  #Initialisations:
22  #*************************************#
23
24  #NumSim denotes the number of data sets to be drawn within one complete
        simulation.
25  NumSim <-10^4
26
27  # T denotes the sample size of the NumSim drawn data sets.
28  T <-60
29
30  #N number of portfolios to be analized.
31  N <- 25
32
33  #Initializing matrix F, it is used to store the T sets of generated
        factors with the mvrnorm function.
34  F <-matrix(0,T,3)
35
36  #Initializing result matirces for gamma 0 (if applicable), 1,2 and 3
37  #The first "c or o" indicates whether the first-pass regression is run
        with constant (c) or without (o).
38  mat.Gamma.mean.cc <- matrix(0, NumSim, 4)
39  mat.Gamma.mean.oc <- matrix(0, NumSim, 4)
40  mat.Gamma.mean.co <- matrix(0, NumSim, 3)
41  mat.Gamma.mean.oo <- matrix(0, NumSim, 3)
42
43  #Matrix to store all betas estimated over the whole observation period.
44  #Referred to as matrix "B" in the text.
45  B = matrix(0,N,3)
46
47  #Matrix where all the residuals of the beta estimation over the whole
```

```
       observation period is stored
48  #Referred to as matrix "E" in the text
49  Res = matrix(0, NROW(Data3x3), N)
50
51  #mat.3F denotes the observed covariance between the 3 factor for the
        whole observation period.
52  mat.cov3F <-matrix (0,3,3)
```

The parameters to be calibrated are:

- True values of $\gamma_1, \gamma_2, \gamma_3$: The true mean of the whole observation period of the three factors (mean of the market portfolio excess return ($Mu\_f$ = true $\gamma_1$), mean of the size factor ($Mu\_smb$= true $\gamma_2$) and the mean of the value factor ($Mu\_hml$ = true $\gamma_3$) as shown in code-lines 58 to 64).

- The covariance of these three factors is evaluated and stored in the matrix $mat.3F$ (code-lines 67,68).

- $B$: The estimation of the betas of the three factors for the whole observation period. This is achieved by running a time-series regression on the return data-series for each portfolio with the observed factors as the explanatory variables. For this purpose, the function $lm()$ is applied within a 1 to N loop to estimate the betas for every single portfolio and to store the results in matrix $B$ (code-line 78). This function slows the whole simulation down substantially compared to applying pure matrix algebra, as it calculates many more results besides the coefficients. However, speed is not of concern for this simulation. One might adjust this if more than 100,000 simulation sets should be drawn. Furthermore, as mentioned in section 3.2, the estimation of the betas without an intercept is done in accordance to the two underlying simulation studies.

- $E$: The just mentioned regressions are also used to derive the residuals of each estimation. Of these residuals the covariance among them is estimated and stored in matrix $E$ as shown in code-line 87.

**Listing 2:** Calibration of Parameters

```
53  #**********************************#
54  #Calibration of the Parameters:
55  #**********************************#
56
57  # Mean of market portfolio return.
58  Mu_f <- mean(Data3x3$Mkt.RF)
59
60  # Mean of value factor.
61  Mu_hml <-mean(Data3x3$HML)
62
```

```
63  # Mean of size factor.
64  Mu_smb <-mean(Data3x3$SMB)
65
66  # Deriving covariance matrix of the three factors Excess return of
        market portfolio, size factor and value factor.
67  mat.3F <- matrix(c(Data3x3$Mkt.RF,Data3x3$SMB,Data3x3$HML),NROW(Data3x3
        ),3)
68  mat.cov3F <- cov(mat.3F)
69
70  #Estimation of the true portfolio betas for all 25 portfolios:
71  #Loop to estimate the "true" beta of the three factors for each
        individual portfolio.
72  #Portfolio number =  i + 1, since the first column of Data5x5 contains
        dates. (eg. 20162)
73  for(i in 1:N) {
74    #Running time-series regression on whole observation period to
          estimate the three betas.
75    lm.result <- lm(Data5x5[,(i+1)]~0 + Data3x3$Mkt.RF + Data3x3$SMB +
          Data3x3$HML)
76
77    #Storing the estimated betas for the N portfolio in matrix B.
78    B[i,] <- c(lm.result$coefficients[[1]], lm.result$coefficients[[2]],
          lm.result$coefficients[[3]])
79
80    #Storing the residuals of the Beta-estimates for the 25 portfolio in
          matrix Res.
81    Res[,i]<-lm.result$residuals
82  }
83
84  # Estimating the covariance of of the residuals across the N Portfolios
        .
85  # Storing the result in matrix "E"
86  # Uncomment "*100" in order to create distorted data sets.
87  E <- cov(Res)#*100
```

Once all the parameters are calibrated, the return data $R$ is generated. In order to generate $T$ returns, $T$ sets of the three factors are required. These can be generated by making use of the *mvrnorm* function available in the library *MASS*[12]. The inputs are the true values of the three factors as well as the covariance among themselves ($mat.3F$) over all 625 observations. The resulting $Tx3$ matrix $F$ (code-line 101) contains $T$ normally distributed sets of the three factors. $F$ is then used to simulate $T$ returns for the 25 portfolios in the following way.

$$FB' + S = R, \tag{32}$$

---

[12]This library is a part of the environment of the R-project programming language.

To illustrate that this return generating process is well defined, the matrix algebra for equation (32) is explained in (33) to (34) by following Wooldridge (2013, p. 802). This algebra corresponds to the code-lines 88 to 108.

$$FB' = \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ \vdots & \vdots & \vdots \\ f_{T1} & f_{T2} & f_{T3} \end{bmatrix} * \begin{bmatrix} \beta_{11} & \dots & \beta_{1N} \\ \beta_{21} & \dots & \beta_{2N} \\ \beta_{31} & \dots & \beta_{3N} \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^{3} f_{ik}\beta_{kn} \end{bmatrix} \tag{33}$$

Where

$i = 1, ..., T$ and $n = 1, ..., N$

$f_{i1} = i^{th}$ generated market factor

$f_{i2} = i^{th}$ generated size factor

$f_{i3} = i^{th}$ generated value factor

$\beta_{1n} = $ observed market beta of $n^{th}$ portfolio

$\beta_{2n} = $ observed size beta of $n^{th}$ portfolio

$\beta_{3n} = $ observed value beta of $n^{th}$ portfolio

$\sum_{k=1}^{3} f_{ik}\beta_{kn}$ is the $(i, n)^{th}$ element of the $TxN$ matrix FB' and equals to the $i^{th}$ simulated return of portfolio $n$ without an error term.

For example, the element $(1, 1)$ of the $FB'$ matrix equals

$$\sum_{k=1}^{3} f_{1k}\beta_{1n} = f_{11}\beta_{11} + f_{12}\beta_{21} + f_{13}\beta_{31}, \tag{34}$$

At this point, the $FB'$ matrix contains $T$ generated portfolio returns for $N$ portfolios. However, so far, the simulated returns do not contain any error terms, which would be an unrealistic assumption for the test statistics since in the underlying data, error terms most likely exist and they are cross-sectionally correlated. Therefore, randomly normally distributed error terms are generated by using the *mvrnorm* function again. The means are 0 and they covary as estimated in matrix $E$. The output of the function is the matrix $TxN$ matrix $S$, which contains the $T$ cross-correlated error terms $e$ for $N$ portfolios. Thus, the final generated returns of matrix $R$ in equation (29) denote as

$$R = \begin{bmatrix} \sum_{k=1}^{3} f_{ik}\beta_{kn} \end{bmatrix} + S = \begin{bmatrix} (\sum_{k=1}^{3} f_{ik}\beta_{kn}) + e_{in}, \end{bmatrix} \tag{35}$$

where all the previous variables are defined as in equation 33 and $e_{in}$ is the $(in)^{th}$ element of the matrix $S$. In other words, $e_{in}$ is the error term for the $i^{th}$ generated return of the $n^{th}$ portfolio.

Since for the second-pass regression, the averages of the $T$ generated returns for all $N$ portfolios are used instead of the $N$ return vectors, $\bar{R}$ is computed (as shown in line 108).

**Listing 3:** Generation of Returns

```
88  #****************************************#
89  #Simulation:
90  #****************************************#
91  #In order to be able to make use of the mvrnrom function the MASS
        library has to be loaded.
92  library(MASS)
93
94  #Running NumSim loops
95  for(k in 1:NumSim) {
96
97    #----------------------------#
98    # Generating Returns:
99
100   # Generating T artificial samples of multivariate normally
          distributed risk-factors, taking into account the observed
          covariance of these three factors among themselves.
101   F <- mvrnorm(T, c(Mu_f,Mu_smb,Mu_hml), mat.cov3F)
102
103   # Generating artificial samples of multivariate normally distributed
          error-terms, taking into account the observed covariance of these
          error terms accros the N portfolios with mean = 0.
104   S <- (mvrnorm(n = T, rep(0,N), E))
105
106   # Generating T artificial returns for N portfolios based on the
          previously generated factors and the true betas and adding
          artificially generated error terms which are correlated accros
          the 25 portfolios.
107   R <- matrix(F%*%t(B)+S,T,N)
108   R.bar <- matrix(apply(R, 2, mean),N,1)
```

Now that $T$ returns with cross-sectional correlation in the residuals are generated for the $N$ portfolios, the first-pass (time-series) regressions are run on the returns in matrix $R$, with and without the intercept restriction. Therefore, for every drawn data set, the betas are estimated in two ways. If the betas are to be estimated without a constant, $F$ takes the form

$$
\begin{bmatrix}
f_{11} & f_{12} & f_{13} \\
f_{21} & f_{22} & f_{23} \\
\vdots & \vdots & \vdots \\
f_{T1} & f_{T2} & f_{T3}
\end{bmatrix}
, \text{ else the form }
\begin{bmatrix}
1 & f_{11} & f_{12} & f_{13} \\
1 & f_{21} & f_{22} & f_{23} \\
\vdots & \vdots & \vdots & \vdots \\
1 & f_{T1} & f_{T2} & f_{T3}
\end{bmatrix}
.
$$

The corresponding adjustment of the $F$- matrix is displayed on code-lines 114 and 120. The betas, as the OLS estimator, of the first-pass regression can be derived in the following way in the case of a multivariate OLS regression (Wooldridge, 2013, p. 802).

$$\hat{\beta} = (F'F)^{-1}F'R \tag{36}$$

The matrices $F$ and $R$ are defined as before. Generally speaking, the resulting $\beta$-estimate matrix takes the following form ($N = 25, K = 3$).

$$\hat{\beta}_o = \begin{bmatrix} \hat{\beta}_{11} & \hat{\beta}_{12} & \hat{\beta}_{13} \\ \hat{\beta}_{21} & \hat{\beta}_{22} & \hat{\beta}_{23} \\ \vdots & \vdots & \vdots \\ \hat{\beta}_{31} & \hat{\beta}_{32} & \hat{\beta}_{NK} \end{bmatrix} or, \hat{\beta}_c = \begin{bmatrix} c_1 & \hat{\beta}_{1,1} & \hat{\beta}_{1,2} & \hat{\beta}_{1,3} \\ c_2 & \hat{\beta}_{2,1} & \hat{\beta}_{2,2} & \hat{\beta}_{2,3} \\ \vdots & \vdots & \vdots & \vdots \\ c_{25} & \hat{\beta}_{25,1} & \hat{\beta}_{3,2} & \hat{\beta}_{NK} \end{bmatrix}.$$

They contain the estimates of the three factor betas for the $N$ portfolios estimated on $T$ simulated returns in $R$ of the corresponding portfolio. $\hat{\beta}_o$ is computed in code-line 116 ($mat.beta.est.o123$) and $\hat{\beta}_c$ is computed in code-line 124 ($mat.beta.est.c123$). Both of the matrices had to be transposed in order to have the right dimensions for the second-pass regression.

**Listing 4:** Estimation of Betas

```
109 #-----------------------------#
110 #Estimating Betas With and Without
111
112 #Estimating betas with an intercept.
113 #Re-assigning F matrix in order to make sure that it has the right
        porperties to run the first-pass OLS regression without an
        intercept.
114 F <-matrix (F,T,3)
115 #Regression on T artificially generated returns R for N portfolios,
        with the T generated sets of the 3 factors stored in matrix F.
116 mat.beta.est.o123 <-t(matrix((solve(t(F)%*%F))%*%t(F)%*%R,3,N))
117
118 #Estimating betas without an intercept
119 #Re-assigning F matrix in order to make sure that it has the right
        porperties to run the first-pass OLS regression without an
        intercept.
120 F <- matrix(c(rep(1,T), F),T,4)
121
122 #Regression on T artificially generated returns R for N portfolios,
        with the T generated sets of the 3 factors + a constant (1) stored
        in matrix F.
123 #Only the coefficients of this regression are stored in mat.beta.est.
        c123, the intercept is not stored.
124 mat.beta.est.c123 <-t(matrix(solve(t(F)%*%F)%*%t(F)%*%R,4,N)[c(2,3,4)
```

```
      ,])
```

Now that the beta estimates have been computed, the cross-sectional regression can be run on the two different sets of beta estimates. Again, the OLS regression are once run with and once without the intercept restriction. Therefore, four final result sets are computed. For this purpose, as done with matrix $F$, the matrices of the explanatory variables have to be adjusted accordingly. The OLS estimators ($\hat{\Gamma}$) are derived by applying the matrix equation described by Wooldridge (2013, p. 802). Thus,

$$\hat{\Gamma} = (C'C)^{-1}C'\bar{R}, \tag{37}$$

where $C$ denotes 25 sets of the three factor betas estimated with a constant. In code-line 131 it is extend on the left with a vertical vector of constants in order to run equation (37) with a constant (method "w/, w/"). Additionally, instead of $C$, $O$ is plugged into (37) to estimate ($\hat{\Gamma}$) with betas which themselves have been estimated without a constant in the first-pass regression. As $C$, $O$ is adjusted in code-line 131 (Method "w/o, w/") and 153 (method "w/, w/") depending on whether ($\hat{\Gamma}$) shall be estimated with or without a constant.

For every method applied, the resulting vector $\hat{\Gamma}$ is stored in the previously initialized result matrices on code-lines 38 to 41. In the case of the methods "w/, w/" and "w/o, w/" $\hat{\Gamma} = \gamma_0, \gamma_1, \gamma_2, \gamma_3$ else $\hat{\Gamma} = \gamma_1, \gamma_2, \gamma_3$. Once the last mean of out of the four variants is stored in these result matrices, the whole simulation, starting with the artificial generation of returns up until the storing of the mean estimates out of the drawn sample is returned $NumSim - 1$ times (e.g. within this thesis the whole process is looped 9,999 times). As soon as these loops are through, the results are plotted as shown in the last snippet of code within this appendix.

**Listing 5:** Estimation of Gammas

```
124 #----------------------------#
125 #Estimating gamma for the four different Methods
126 #The matrices C and O are used to run the variants of these second-pass
        regressions on R.
127 #C indicates that the betas in the matrix of were measured with a
        constant.
128 #O indicates that the betas in the matrix of were measured without a
        constant
129 #G, a TxN matrix, contains the resulting Gamma estimates for the N
        portfolios.
130
131 #Method w/, w/
132 #Estimating gamma 0,1,2,3 out of estimated betas with an intercept
133 C <- ( matrix ( c ( rep ( 1 ,N) , mat . beta . est . c123 ) ,N,4 ) )
134 #Run the regression on R.bar with C as the explantory.
135 G <-solve ( t (C)%*%C)%*%t (C)%*% R. bar
136 #Store the mean of the T estimated gammas in the corresponding result
        vector.
```

```
137 mat.Gamma.mean.cc[k,] <- G
138
139 #Method w/o, w/
140 #Estimating gamma 0,1,2,3 out of estimated betas without an intercept
141 O <- (matrix(c(rep(1,N),mat.beta.est.o123),N,4))
142 #Run the regression on R.bar with O as the explantory.
143 G <-solve(t(O)%*%O)%*%t(O)%*% R.bar
144 #Store the mean of the T estimated gammas in the corresponding result
        vector.
145 mat.Gamma.mean.oc[k,] <- G
146
147 #Method w/, w/o
148 #Only estimating gamma 1,2,3 out of estimated betas with an intercept
149 C <- (mat.beta.est.c123)
150 #Run the regression on R.bar with C as the explantory.
151 G <- solve(t(C)%*%C)%*%t(C)%*% R.bar
152 #Store the mean of the T estimated gammas in the corresponding result
        vector.
153 mat.Gamma.mean.co[k,] <- G
154
155 #Method w/, w/
156 #Only estimating gamma 1,2,3 out of estimated betas without an
        intercept
157 O <- (mat.beta.est.o123)
158 #Run the regression on O.bar with C as the explantory.
159 G <- solve(t(O)%*%O)%*%t(O)%*% R.bar
160 #Store the mean of the T estimated gammas in the corresponding result
        vector.
161 mat.Gamma.mean.oo[k,] <- G
162
163 #Repeat Loop NumSim-1 times. (e.g. 9,999 times)
164 }
```

For the code to display the results, the comments are only made in-code, as it is not part of the simulation. In general, the result matrices in code-lines 38 to 41 contain $NumSim$ (=10,000) estimates for 3 or 4 gammas (depending on the method). The results are displayed through $boxplots$. Each $boxplot$ represents one column in the four result matrices. All the graphs plotted are automatically stored in the working directory of the user. The graphs do not automatically show in the plot field of $R$-Studio.

**Listing 6:** Display of Results

```
165 #***********************************#
166 #Display of Results:
167 #***********************************#
168
```

```
169 # The graphs are automatically stored as a PDF in the working directory
        of the code.
170 # They are named by their premium + sample size
171 # The graphs are not automatically shown in R-Studio
172 # To display a graph in R-Studio call dev.new() and do not call dev.off
        ()
173 h <- 5 #set hight of the to be stored graphs
174 w <- 3.5 #set width of the to be stored graphs
175 NameVar <- "distorted" #Choose a caption in order to give the to be
        named graph a unique name.
176 options(scipen=10)
177
178 #Create Boxplot for Gamma 0 for the cases "w/ w/" and "w/ow" and store
        it in the working directory.
179 name <- paste(c("gamma0_",NameVar, T, ".pdf"), collapse = "")
180 pdf(name, width = w, height= h)
181 boxplot(mat.Gamma.mean.cc[,1], mat.Gamma.mean.oc[,1],
182         ylab="",
183         ylim=c(-4,4),
184         las=2,
185         mgp=c(2,1,.5),
186         names=c("w/, w/", "w/o w/"))
187 mtext(paste("T =", T, ", NumSim =", NumSim))
188 abline(h=0.0, col = "black", lty = "dashed", lwd=1.5)
189 title(ylab=bquote(gamma[0] ~ "(in %)"), line=1.55, cex.lab=1.2)
190 dev.off()
191
192 #Create Boxplot for Gamma 1 and store it in the working directory.
193 name <- paste(c("gamma1_",NameVar, T, ".pdf"), collapse = "")
194 pdf(name, width = w, height= h)
195 boxplot(mat.Gamma.mean.cc[,2], mat.Gamma.mean.oc[,2], mat.Gamma.mean.co
        [,1],mat.Gamma.mean.oo[,1],
196         ylab="",
197         ylim=c(-4,4),
198         las=2,
199         names=c("w/, w/", "w/o, w/ ", "w/, w/o ","w/o, w/o "))
200 mtext(paste("T =", T, ", NumSim =", NumSim))
201 abline(h=Mu_f, col = "red", lty = "dashed", lwd=1.5)
202 abline(h=0.0, col = "black", lty = "dashed", lwd=1.5)
203 title(ylab=bquote(gamma[1] ~ "(in %)"), line=1.55, cex.lab=1.2)
204 axis(2.0, at=c(round(Mu_f,2)), las=1, labels=FALSE, line = -0.75)
205 axis(2.0, at=(round(Mu_f,2)), las=1, labels=TRUE, line = -0.75 )
206 dev.off()
207
208 #Create Boxplot for Gamma 2 and store it in the working directory.
```

```
209  name <- paste(c("gamma2_",NameVar, T, ".pdf"), collapse = "")
210  pdf(name, width = w, height= h)
211  boxplot(mat.Gamma.mean.cc[,3], mat.Gamma.mean.oc[,3], mat.Gamma.mean.co
         [,2], mat.Gamma.mean.oo[,2],
212          ylab="",
213          ylim=c(-4,4),
214          las=2,
215          names=c("w/, w/", "w/o, w/ ", "w/, w/o ","w/o, w/o "))
216  mtext(paste("T =", T, ", NumSim =", NumSim))
217  abline(h=Mu_smb, col = "green", lty = "dashed", lwd=1.5)
218  abline(h=0.0, col = "black", lty = "dashed", lwd=1.5)
219  title(ylab=bquote(gamma[2] ~ "(in %)"), line=1.55, cex.lab=1.2)
220  axis(2.0, at=c(round(Mu_smb,2)), las=1, labels=FALSE, line = -0.75)
221  axis(2.0, at=(round(Mu_smb,2)), las=1, labels=TRUE, line = -0.75)
222  dev.off()
223
224  #Create Boxplot for Gamma 3 and store it in the working directory.
225  name <- paste(c("gamma3_",NameVar, T, ".pdf"), collapse = "")
226  pdf(name, width = w, height= h)
227  boxplot(mat.Gamma.mean.cc[,4], mat.Gamma.mean.oc[,4], mat.Gamma.mean.co
         [,3], mat.Gamma.mean.oo[,3],
228          ylab="",
229          ylim=c(-4,4),
230          las=2,
231          names=c("w/, w/", "w/o, w/ ", "w/, w/o ","w/o, w/o "))
232  mtext(paste("T =", T, ", NumSim =", NumSim))
233  abline(h=Mu_hml, col = "green", lty = "dashed", lwd=1.5)
234  abline(h=0.0, col = "black", lty = "dashed", lwd=1.5)
235  title(ylab=bquote(gamma[3] ~ "(in %)"), line=1.55, cex.lab=1.2)
236  axis(2.0, at=c(round(Mu_hml,2)), las=1, labels=FALSE, line = -0.75)
237  axis(2.0, at=(round(Mu_hml,2)), las=1, labels=TRUE, line = -0.75)
238  dev.off()
239
240
241  #Plot beta distribution of the 25 portfolios + the average of the
         montly returns of the 25 portfolios.
242  BR <- cbind(B, (apply(Data5x5[,2:26],2,mean)))
243  pdf("beta.pdf", width = 9)
244  matplot(BR,type =c("b"),lty=c(1:2,4:5),pch=1:4,col=1,
245          ylab = "Beta",
246          xlab = "Portfolio")
247  legend("topright",legend = c(expression(beta[M]), expression(beta[SMB])
         , expression(beta[HML]), expression(R[P])), col=1, pch=1:4)
248  dev.off()
249
```

```
250  #Output of table 1 data
251  apply ( mat.Gamma.mean.oo ,2 ,mean )
252  apply ( mat.Gamma.mean.oo ,2 ,sd )
253
254  apply ( mat.Gamma.mean.co ,2 ,mean )
255  apply ( mat.Gamma.mean.co ,2 ,sd )
256
257  apply ( mat.Gamma.mean.oc ,2 ,mean )
258  apply ( mat.Gamma.mean.oc ,2 ,sd )
259
260  apply ( mat.Gamma.mean.cc ,2 ,mean )
261  apply ( mat.Gamma.mean.cc ,2 ,sd )
```