

Essays on Risk and Return in the Stock Market

by

Olivier Richard Henri Ledoit

Submitted to the Sloan School of Management
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy in Management

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 1995

© Massachusetts Institute of Technology 1995. All rights reserved.

Author
Sloan School of Management
June 1, 1995

Certified by
Andrew W. Lo
Harris & Harris Group Professor of Finance
Thesis Supervisor

Accepted by
Birger Wernerfelt
Chairman, Departmental Committee on Graduate Students

ARCHIVES

MASSACHUSETTS INSTITUTE
OF TECHNOLOGY

JUN 08 1995

Essays on Risk and Return in the Stock Market

by

Olivier Richard Henri Ledoit

Submitted to the Sloan School of Management
on June 1, 1995, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy in Management

Abstract

The theory of risk and return in the stock market is perhaps the best understood case of economic decision under uncertainty. The prominent developments in this field are the procedure for portfolio selection, the Capital Asset Pricing Model (CAPM), and the Arbitrage Pricing Theory (APT). This dissertation revisits all three of these important themes, with special attention paid to developing new angles of attack.

The first essay proposes an original solution to a troublesome statistical problem often encountered when implementing portfolio selection: the estimation of the inverse of the variance-covariance matrix with less observations than variables; an empirical application of this new statistical technique challenges recent claims of flatness in the relationship between expected stock returns and CAPM betas. The second chapter is a theoretical essay introducing a strengthened version of the APT that expresses rigorously the APT's ideas when the number of stocks is finite. This theoretical framework is used in the third essay to compute an upper bound on deviations from beta pricing, both for the CAPM and the APT; empirical results show that this bound is rather loose, indicating that betas can easily fail to characterize expected returns to any acceptable degree of accuracy; furthermore, when estimation error on risk premia is acknowledged, the optimal number of factors in the APT is very small (three or less), and the error is roughly as large in the APT's beta pricing equation as in the CAPM's.

Thesis Supervisor: Andrew W. Lo

Title: Harris & Harris Group Professor of Finance

Acknowledgments

It is quite an exaggeration to write down my sole name as the author of this thesis. Sometimes I feel like all I did was channel the energies of all the people who believed in me. Together, they probably sacrificed as much as I did to this dissertation. I want to thank them here. I thank all the people who have helped me, and in particular:

My mother, for raising me all on her own, loving me, teaching me that school is important, and letting me move far away from her;

My grandmother, for loving me and being so wise;

Sally Whitman, for holding my hand all the way through;

Tim Crack, for helping me keep my sanity, forcing me to be rigorous, teaching me how to write, showing me how to teach, giving feedback on this research, but most important of all for becoming my friend over the course of these four hard years;

Jean-Luc Vila, for helping me get into the program, making me feel welcome at MIT, transmitting his enthusiasm for Finance, and offering a job at Goldman Sachs that I refused;

Andrew Lo, for accepting to be my advisor, figuring out how to handle me, greatly contributing to the quality of the first chapter, and helping me get the job I wanted;

John Heaton, for giving a lot of very thoughtful feedback, and supervising the last chapters;

Bin Zhou, for contributing significantly to the first chapter, and being on my committee;

John Cox, for sharing some of his experience;

Jack Silverstein, for explaining the spectral theory of large random matrices, which inspired the resolution of the central problem of the first chapter;

Dimitri Bertsimas and Whitney Newey, for detailed comments;

George Constantinides, Phil Dybvig, Bob Litzenberger, Craig MacKinlay, Dick Roll, Steve Ross and Rob Stambaugh, for illuminating conversations;

Petr Adamek, for always accepting the challenge of a solid technical discussion;

Angel Serrat Tubert, for being of good advice on non-technical, yet important, matters;

Firooz Partovi, the Laboratory for Financial Engineering, the Sloan Computer Lab, and
Athena, for outstanding computing support;

Mary Marshall and Sharon Cayley, for always being on my side instead of the administration's;

The whole Fall 1994 MIT Finance faculty, for valuable comments on my thesis proposal;

My fellow Finance Ph.D. students at MIT, for valuable comments also;

My future colleagues at UCLA, for making my last six months as a Ph.D. student much
happier, and I am sure my following years as Assistant Professor too!

Participants to seminars at the NBER, MIT's Finance, Economics and Operations Research departments, UCLA, Washington University in Saint Louis, Yale, Chicago and Wharton, for useful discussions;

And, last but not least, the area's support staff, for being nice to me.

It is a long list, but it is not nearly long enough, so I apologize to the people not cited above.
All your support was appreciated, is gratefully acknowledged, and will be remembered.

I dedicate this thesis to my mother, a smart woman who never finished high school, and
always wanted me to do better. This is her accomplishment as much as my own.

Contents

1	Portfolio Selection: Improved Covariance Matrix Estimation	11
1.1	Introduction	12
1.1.1	Overview	12
1.1.2	Comparison with Existing Literature	14
1.2	Sample Covariance Matrix	15
1.2.1	Model	16
1.2.2	Norm	17
1.2.3	Consistency	19
1.2.4	Portfolio Selection and Covariance Matrix Eigenvalues	20
1.2.5	Sample Covariance Matrix Eigenvalues	22
1.2.6	Particular Case: the Identity Matrix	23
1.2.7	Potential Applications to Tests for the Number of Factors in the APT	25
1.3	Improved Covariance Matrix Estimation	26
1.3.1	Linear Shrinkage of Sample Eigenvalues	26
1.3.2	Optimal Linear Shrinkage	27
1.3.3	Generalization	30
1.3.4	Comparison with Previous Work in Multivariate Statistics	32
1.4	Application to Portfolio Selection	34
1.4.1	Monte-Carlo Simulations	34
1.4.2	Historical Data	38
1.4.3	Testing an Implication of the CAPM	41
1.5	Conclusion	44

2	The δ-Arbitrage Pricing Theory	46
2.1	Introduction	46
2.2	Arbitrage Pricing Theory	48
2.2.1	Review of the APT	48
2.2.2	Limitation of the APT	49
2.3	δ -Arbitrage	50
2.3.1	Definitions of Arbitrage	50
2.3.2	δ -Arbitrage	51
2.3.3	Economic Justification	52
2.3.4	Non-Marketable Assets	54
2.3.5	Empirical Evidence	54
2.3.6	Choice of δ	55
2.4	δ -Arbitrage Pricing Theory	56
2.4.1	Formulation	56
2.4.2	Relation to Existing Asset Pricing Theories	58
2.4.3	Economic Contents	59
2.4.4	Testability	60
2.5	Choice of Factors	61
2.5.1	Factors vs. Residuals	61
2.5.2	Exogenous Factors	62
2.5.3	Covariance Matrix Eigenvectors	63
2.6	Estimation Error	63
2.6.1	Forecasting Residual Space	64
2.6.2	Maximum Residual Eigenvalue	64
2.6.3	Optimal Number of Factors	65
2.7	Conclusion	67
3	Is Beta Pricing Accurate?	69
3.1	Empirical Strategy	70
3.1.1	Beta Pricing	70

3.1.2	Objective	72
3.1.3	Sources of Deviation from Beta Pricing	73
3.1.4	Data	74
3.2	Exogenous Market Factor	74
3.2.1	Factor vs. Residuals	75
3.2.2	Residual Risk	75
3.2.3	Risk Premium Estimation Error	76
3.2.4	Overall Evaluation	77
3.3	Eigenvectors	77
3.3.1	Factors vs. Residuals	77
3.3.2	Residual Risk	78
3.3.3	Risk Premium Estimation Error	78
3.3.4	Overall Evaluation	79
3.4	Conclusion	79
A	Spectral Theory of Large Random Matrices	80
A.1	Mathematical Tools	80
A.2	Asymptotic Results	82
A.3	From True to Sample Eigenvalues	84
A.4	From Sample to True Eigenvalues	84
B	Other Structured Estimators	87
B.1	All Variances, Respectively Covariances, Are Equal	87
B.2	Diagonal Matrix	87
B.3	All Correlation Coefficients Are Equal	88
B.4	Single Index Model	88
C	Proofs of Chapter 1	89
C.1	Theorem 1	89
C.2	Theorem 2	90
C.3	Theorem 3	92

C.4	Theorem 4	92
C.5	Theorem 5	93
C.6	Theorem 6	93
C.7	Theorem 7	99
C.8	Theorem 8	102
C.9	Theorem 9	102
C.9.1	$\ \hat{\hat{\Sigma}} - \hat{\Sigma}\ ^2 \xrightarrow{P} 0$	102
C.9.2	$E[\ \hat{\hat{\Sigma}} - \Sigma\ ^2] - E[\ \hat{\Sigma} - \Sigma\ ^2] \rightarrow 0$	103
C.9.3	$(\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2) - (r_1^2 r_2^2 / d^2) \xrightarrow{P} 0$	104
C.10	Theorem 10	105
D	Proofs of Chapter 2	106
D.1	Theorem 11	106
D.2	Theorem 12	107
D.3	Theorem 13	107
D.4	Theorem 14	108
D.5	Theorem 15	109
D.6	Theorem 16	110
E	Tables	111
F	Figures	116
	Bibliography	131

List of Tables

1.1	Result of 1,000 Monte-Carlo Simulations for Central Parameter Values. . .	36
E.1	Comparison of the Ex-Post Standard Deviations of Ex-Ante Minimum Variance Portfolios.	111
E.2	Comparison of the Ex-Post Standard Deviations of Minimum Variance Portfolios.	112
E.3	Predictive OLS Cross-Sectional Regression of Returns on Betas over 1936- 1992.	113
E.4	Predictive GLS Cross-Sectional Regression of Returns on Betas over 1936- 1992.	114
E.5	Predictive GLS Cross-Sectional Regression of Returns on Betas over 1936- 1992.	115

List of Figures

F-1	Sample vs. True Eigenvalues.	117
F-2	Geometric Interpretation of Theorem 5.	118
F-3	Bayesian Interpretation.	119
F-4	Effect of the Ratio of Number of Variables to Number of Observations on the Percentage Relative Improvement in Average Loss (PRIAL).	120
F-5	Effect of the Dispersion of Eigenvalues on the Percentage Relative Improvement in Average Loss (PRIAL).	121
F-6	Effect of the Product of Variables by Observations on the Percentage Relative Improvement in Average Loss (PRIAL).	122
F-7	Weights on Structured Estimators.	123
F-8	Ex-Post Characteristics of Ex-Ante Constrained Minimum Variance Portfolios.	124
F-9	Domain where the Value of $s_{L\hat{H}}$ is Known from Equation (A.4).	125
F-10	Top 100 Eigenvalues of the Covariance Matrix of Stock Returns.	126
F-11	Ratio of Consecutive Eigenvalues of the Covariance Matrix of Stock Returns.	127
F-12	Ratio of the Top Eigenvalue to Lesser Eigenvalues.	128
F-13	Beta Pricing Error Bound.	129
F-14	Accuracy vs. Parsimony.	130

Chapter 1

Portfolio Selection: Improved Covariance Matrix Estimation

This chapter studies the estimation of the covariance matrix of returns on all stocks traded in the stock market, for portfolio selection. The number of observations is assumed to go to infinity, but the standard asymptotic assumption that keeps the number of variables bounded is lifted. In practice, this is appropriate when the number of traded stocks is at least of the same order of magnitude as the number of time periods, which is the usual case.

The first part characterizes intuitively and analytically the behavior of the sample covariance matrix in this case. Some of this work is potentially applicable to tests for the number of factors in the Arbitrage Pricing Theory (APT). The second part develops a simple and versatile estimator that has lower mean squared error than the sample covariance matrix. This estimator provides attractive answers to some fundamental questions in multivariate statistics. In the third and last part, Monte-Carlo simulations and historical data indicate that the new estimator improves over existing ones for portfolio selection: it yields portfolios with significantly lower risk than was previously possible. One of the empirical applications can be interpreted as a test of the Capital Asset Pricing Model (CAPM) with higher power than existing tests. It finds a significant and robust positive relationship between returns and betas, in contrast with less powerful tests in the literature.

1.1 Introduction

1.1.1 Overview

The objective of this study is to estimate the covariance matrix of returns on all stocks traded in the stock market. This is important because the covariance matrix is a necessary input to Markowitz (1952) portfolio selection, a central method in stock market finance.

Our original approach is to assume that the number of observations T goes to infinity, as in standard asymptotics, but relax the standard asymptotic assumption that the number of variables N remains bounded by a constant: we only assume that N is bounded by a constant times T . It is a more realistic approximation of actual stock returns data, because typically the number of traded stocks N is at least of the same order of magnitude as the number of time periods T .

In the first part, we show that the sample covariance matrix is no longer consistent in this framework. Its mean squared error is of order N/T . For example, the sample covariance matrix of $N = 1,000$ stocks based on $T = 2,000$ observations is approximately as erroneous as the variance of the return on $N = 1$ stock estimated from $T = 2$ observations. Not only is the error substantial, but its nature is particularly damaging to portfolio selection: it causes the sample covariance matrix to be near-singular or singular. When the sample covariance matrix is near-singular, inverting it amplifies error and yields grossly inaccurate results for portfolio selection. This is the case if N is of the same order of magnitude as T . When the sample covariance matrix is singular, it cannot be inverted and cannot be used for portfolio selection at all. This is the case if N exceeds T .

We also review the spectral theory of large-dimensional random matrices. This theory gives the relationship between the eigenvalues of true and sample covariance matrices as a function of the ratio N/T , when T goes to infinity. It is the fact that the smallest sample covariance matrix eigenvalues are biased down towards zero that causes the singularity problem. This theory can potentially be used to test hypotheses about the eigenvalues of the covariance matrix of stock returns, such as the ones made by the Arbitrage Pricing Theory (APT).

In the second part, we improve over the sample covariance matrix. Some authors

impose parsimonious structure (e.g. all pairs of stocks have the same correlation coefficient) to obtain an estimator with fewer free parameters. Better yet, Frost and Savarino (1986) combine such a “structured” estimator with the sample covariance matrix. We focus on weighted averages of a structured estimator with the sample covariance matrix and ask: what are the optimal weights? In our asymptotic framework, simple estimators of the weights minimizing mean squared error are consistent. We thus show how to improve both on any given structured estimator and on the sample covariance matrix by combining them in an asymptotically optimal way. Not only does it reduce mean squared error, but it generally escapes the singularity problem.

This method can be interpreted in Bayesian terms. The structured estimator can be called the prior, and its combination with the sample covariance matrix the posterior. Fundamental Bayesian questions have always been: Where does the prior come from? How confident are we in the prior? In finite sample, it is very hard to answer these questions satisfactorily. By contrast, in our asymptotic framework, the prior can be taken as any structured estimator, and the degree of confidence in the prior can be estimated consistently.

In the third part, we show that the improved estimator performs well in practice. In Monte-Carlo simulations, it has lower mean squared error than the sample covariance matrix, even in very small sample. Historical simulations confirm that, for a given set of constraints, our estimator yields portfolios with significantly lower risk than existing estimators.

One of our historical simulations is the first predictive Generalized Least Squares (GLS) cross-sectional regression of stock returns on betas and size. Similar regressions have been interpreted as tests of the CAPM. Thanks to our improved covariance matrix estimator, our GLS-based tests have more power than the tests in the literature, which are based on Ordinary Least Squares (OLS). By contrast with OLS tests, our GLS tests find a significant and robust positive relationship between returns and betas.

In this section, we present an overview of the chapter and contrast it with the existing literature. In Section 1.2, we study the behavior of the sample covariance matrix when the number of variables is allowed to grow large. We develop a family of estimators

that improve over the sample covariance matrix in Section 1.3. In Section 1.4, we see how these estimators perform for portfolio selection. Section 1.5 concludes. Appendix A contains details about the spectral theory of large-dimensional random matrices. Appendix B contains formulas for the more complicated versions of our estimator. Proofs are in Appendix C.

1.1.2 Comparison with Existing Literature

Jobson and Korkie (1980) show that using the sample covariance matrix for portfolio selection can cause severe problems. In some cases, it is better to use the identity matrix instead. Our main intuition is that a well-chosen linear combination of the sample covariance matrix with the identity can work even better than either. Our main contribution is to show how to choose this linear combination well.

Bawa, Brown and Klein (1979) argue that estimation risk coming from sample covariance matrix error is of the same nature as investment risk coming from stock return volatility. Their idea is of a Bayesian nature. One of their recommendations is to combine the sample covariance matrix with an “informative” prior. The more confident we are in the prior, the heavier it should weigh in the combination. They do not show how to obtain the prior and the degree of confidence in it. This is what we do.

Our work is closest in spirit to Frost and Savarino’s (1986). The difference is that they work in finite sample, while we work asymptotically. In finite sample, they have to ignore dependence between the prior and the sample covariance matrix, assume normality, and require that observations outnumber variables. Their formula is not explicit and is costly to compute for large universes of stocks. Asymptotically, we avoid all these problems. The price to pay is that peak performance only kicks in when N and T are large (larger than, say, 30), but this is almost always the case in practice.

Kandel and Stambaugh (1994) analyze cross-sectional regressions of stock returns on betas. The CAPM implies a positive slope. A problem arises because the market, with respect to which betas are measured, is only known approximately (Roll, 1977). Then the regression method matters. With Ordinary Least Squares (OLS), the regression slope can be

anything, even if the CAPM holds. OLS uses the identity in place of the covariance matrix of stock return residuals. With Generalized Least Squares (GLS), however, the estimated regression slope must be close to the one implied by the CAPM, if the CAPM holds and the market proxy is close to the true market. GLS require an estimator of the covariance matrix of residuals.¹ Where to find it? Usually, the sample covariance matrix is out of the question because it is near-singular or singular. We show that a linear combination of the identity and the sample covariance matrix can be used to run GLS regressions.

Brown (1989) finds that APT tests based on sample covariance matrix eigenvalues are extremely sensitive to the relative magnitudes of the number of time periods T and the number of stocks N . His results are obtained by Monte-Carlo simulations in a stylized case. We review an equation that gives the distribution of sample eigenvalues as a function of the distribution of true eigenvalues and the ratio N/T , when T goes to infinity. Potentially, it could be used to correct APT tests for the effect noticed by Brown.

To the best of our knowledge, the only published results on the sample covariance matrix when N goes to infinity with T characterize eigenvalues. This literature is part of the spectral theory of large-dimensional random matrices. Marčenko and Pastur (1967) first obtained its central equation, which is the one that we alluded to in the previous paragraph. The most recent and general result is by Silverstein (1994). We could only find two statistical applications in this literature: Wachter (1976) and Silverstein and Combettes (1992). Both are restricted to special cases, and study only eigenvalues. By contrast, we work in the general case, and are interested in the whole sample covariance matrix.

1.2 Sample Covariance Matrix

We analyze the behavior of the sample covariance matrix when the number of variables is large, the typical case for portfolio selection with stocks.

¹The term GLS sometimes means using the true covariance matrix; here, just an estimator.

1.2.1 Model

Consider a very simple situation where we relax the standard asymptotic assumption that keeps the number of variables fixed.

Assumption 1 *Let $T = 1, 2, \dots$ index a sequence of statistical models. For every T , X_T is an $N_T \times T$ matrix of T independent and identically distributed (iid) observations on a system of N_T random variables with mean zero and $N_T \times N_T$ covariance matrix $\Sigma_T = E[(1/T)X_T X_T']$, where $E[\cdot]$ denotes expectation and prime denotes transposition. The sample covariance matrix is $\tilde{\Sigma}_T = (1/T)X_T X_T'$. Assume that there exists a constant A independent of T such that $N_T \leq A T$.*

All the quantities in this chapter depend on T unless otherwise specified. For fluidity we omit the subscript T . Assumption 1 prevents the number of variables N from growing infinitely faster than the number of observations T .

The assumption that the random variables have mean zero is not restrictive because, in practice, we can always subtract some estimator of mean returns. How to estimate mean returns is strictly outside the scope of this chapter.

Decompose the covariance matrix into eigenvectors and eigenvalues: $\Sigma = U \Lambda U'$, where U is a rotation matrix ($U'U = UU' = I$ the identity matrix) whose columns are the eigenvectors of Σ , and Λ a diagonal matrix whose diagonal elements are the eigenvalues of Σ . Define $Y = U'X$, an $N \times T$ matrix of T iid observations on a system of N uncorrelated random variables that spans the same space as the original system.

We must impose some cross-sectional restrictions in order to obtain results when we allow N to grow without bounds.

Assumption 2 *Let $(y_{11}, \dots, y_{N1})'$ denote the first column of the matrix Y' . The average eighth moment is bounded in the following sense: there exists a constant B independent of T such that $E[(1/N) \sum_{i=1}^N y_{i1}^8] \leq B$.*

Assumption 3 $\text{Cov}[y_{i1}y_{j1}, y_{k1}y_{l1}] = 0$ when the set $\{i, j\}$ does not intersect with the set $\{k, l\}$.

Assumptions 1-3 are implicit throughout the remainder of the chapter.

1.2.2 Norm

The originality of this framework is that the dimension N of the covariance matrix can change as T goes to infinity, and can even go to infinity itself: the space where the covariance matrix lives is changing. This makes the definition of a norm on covariance matrices delicate, but not impossible.

Two solutions come to mind: either define a norm on an infinite-dimensional space into which every finite-dimensional space can be embedded, or define a sequence of norms directly on the finite-dimensional spaces. Since it is not exactly clear how to implement the first solution, I opt for the second one.

The sequence of norms (one norm corresponding to each dimension N) is built around the Frobenius norm, which is often used in linear algebra.

Definition 1 *The norm of the $N \times N$ symmetric matrix S with entries $(s_{ij})_{i,j=1,\dots,N}$ and eigenvalues $(l_i)_{i=1,\dots,N}$ is defined by:*

$$\|S\|^2 = c_N \text{tr}(S^2) = c_N \sum_{i=1}^N \sum_{j=1}^N s_{ij}^2 = c_N \sum_{i=1}^N l_i^2, \quad (1.1)$$

where tr denotes the trace and c_N is a scalar coefficient. This norm is a quadratic form on the linear space of $N \times N$ symmetric matrices. Its associated inner product is: $S_1 \circ S_2 = c_N \text{tr}(S_1 S_2)$, where S_1 and S_2 are $N \times N$ symmetric matrices.

It is attractive for the squared norm of a matrix to accumulate the squares of individual entries. The coefficient c_N controls the asymptotic behavior of the sequence of norms. Rigorously speaking, the symbol for the norm $\|\cdot\|$ should be bearing the subscript N .

In order to complete the construction of the sequence of norms, we must choose what asymptotic properties we want to impose on it, and determine the sequence of coefficients c_N accordingly. Remember that an N -dimensional matrix represents a linear operator on the space of N -dimensional vectors. A desirable property is that the norm of familiar linear operators remains well-behaved as N goes to infinity.

The standard definition of the Frobenius norm uses $c_N = 1$. This may be appropriate for the standard case where the dimension N is fixed, but it would cause severe paradoxes

as N goes to infinity. For example, it would make the norm of the identity matrix go to infinity with N . This is not acceptable because, as a linear operator, the identity leaves vectors unchanged, and this operation is too mild to deserve an infinite norm.

The same paradox can be rewritten in terms of the distance between two matrices, which is equal to the norm of their difference. This is advantageous because $\|\cdot\|$ is only used below as a distance, not a norm.

The problem with $c_N = 1$ is that the distance between two sequence of matrices could increase just because their dimension increases. All other things equal, distances would be greater, the greater the dimension. In a general sense, distances would be larger between two high-dimensional matrices than between two low-dimensional ones. To present an analogy, it would be as ill-advised as measuring in the same unit the distance between two cities and the distance between two galaxies.

This paradox is resolved by defining a *relative* distance. The distance between two N -dimensional matrices is divided by the distance between two benchmark matrices of the same dimension N . Relative distance corrects for the potentially disturbing impact of dimension. The benchmark must be chosen carefully. I take the benchmark as the distance from the null matrix to the identity. This convention determines c_N uniquely.

Definition 2 *The scalar coefficient not specified by Definition 1 is: $c_N = 1/N$.*

Any choice of c_N such that the distance from the identity to the null matrix remains bounded away from zero and infinity would induce a norm equivalent to Definition 2's. This is a very large class, and arguably it contains any distance that would make sense in this context. Equivalence means that the notions of convergence and consistency are blind to the particular distance in the class. We can thus be confident that Definitions 1-2 capture an intuitively satisfying notion of distance.

A simple example illustrates the asymptotic behavior of the distance defined above. Let M_1 denote the $N \times N$ matrix with one in its top left entry and zeros everywhere else. Let M_0 denote the $N \times N$ matrix with zeros everywhere (i.e. the null matrix). M_1 and M_0 differ in a way that is independent of N : the top left entry is not the same. Yet the squared distance $\|M_1 - M_0\|^2 = 1/N$ depends on N .

This apparently surprising remark has an intuitive explanation. M_1 and M_0 disagree on the first dimension, but they agree on the $N - 1$ others. The importance of their disagreement is relative to the extent of their agreement. If $N = 1$, then M_1 and M_0 have nothing in common, and their distance is 1. If $N \rightarrow \infty$, then M_1 and M_0 have almost everything in common, and their distance goes to 0. Thus, disagreeing on one entry can either be important (if this entry is the only one) or negligible (if this entry is lost among many others).

It was important to take the time to define the “right” distance because results about consistency are only as interesting as the distances that they are obtained under. If we want the appealing features of the Frobenius norm, it seems that the above choice is the only one (up to equivalence) that makes any sense as N goes to infinity.

Even though Definition 2 is crucial for theoretical results of consistency, it does not matter at all in practice. As will be seen later, the usefulness of this chapter from an empirical point of view is to estimate consistently shrinkage intensities (the scalars m and r_2^2/d^2 , see Section 1.3.2) that are *ratios* of distances or inner products of N -dimensional matrices. Therefore the scalar coefficient c_N will cancel itself out from every formula used in practice.

1.2.3 Consistency

Let $m = \Sigma \circ I$, where I is the identity. The scalar m measures the scale of the covariance matrix. m is the average of the diagonal elements and also the average of the eigenvalues of Σ . The scalar multiple of the identity closest to Σ is mI . mI is the orthogonal projection of Σ onto the line spanned by I . If $I \circ I = \|I\|^2$ was not equal to one, then the correct definition would be: $m = (\Sigma \circ I)/(I \circ I)$.

The mean squared error of the sample covariance matrix is of order N/T .

Theorem 1 $E[\|\tilde{\Sigma} - \Sigma\|^2] = (N/T) m^2 \rightarrow 0$, where convergence is meant as T goes to infinity.

When N/T does not vanish, which is the general case under Assumption 1, the sample covariance matrix is not consistent. When N/T vanishes, which is a special case of

Assumption 1, the sample covariance matrix is consistent. In particular, when N is bounded, our framework degenerates to standard asymptotics.

$\tilde{\Sigma}$ is not consistent because of its off-diagonal elements. Granted, the variance of each one of them vanishes in $1/T$, but so many of them accumulate that the error of $\tilde{\Sigma}$ as a whole does not vanish.

$T = 2,000$ time periods might sound like a lot, but it is not enough if we have as many as $N = 1,000$ stocks: it is about as bad as using two observations to estimate the variance of one random variable. 1,000 is less than half the number of stocks trading on the New York Stock Exchange (NYSE) alone. In order to estimate a $1,000 \times 1,000$ covariance matrix accurately, we need at least, say, 10,000 observations, which means 40 years of daily data, longer than the Center for Research in Security Prices (CRSP) database holds, and in any case long enough for nonstationarity to become a major concern.

Even though we have not tried to obtain a formal proof, we firmly believe that no other covariance matrix estimator is consistent under Assumptions 1-3. Yet all hope is not lost. More than its existence, it is the nature of this error that hurts portfolio selection. We will soon see that the heart of the problem lies in the smallest eigenvalues of the sample covariance matrix. First, we review the importance of covariance matrix eigenvalues for portfolio selection.

1.2.4 Portfolio Selection and Covariance Matrix Eigenvalues

Markowitz (1952) considers the problem of selecting the $N \times 1$ vector of weights w of a portfolio of N stocks whose returns have $N \times N$ covariance matrix Σ , under the K linear constraints defined by the $N \times K$ matrix of coefficients C and the $K \times 1$ right-hand-side vector γ . The objective is to minimize the variance of portfolio returns:

$$\begin{aligned} \min_w & w' \Sigma w \\ \text{s.t.} & C' w = \gamma \end{aligned} \tag{1.2}$$

$$\rightarrow w = \Sigma^{-1} C (C' \Sigma^{-1} C)^{-1} \gamma \tag{1.3}$$

Typical constraints impose that weights sum to one and portfolio returns have a required expectation.

Recall the decomposition $\Sigma = U\Lambda U'$. Let u_1, \dots, u_N denote the columns of U , i.e. the eigenvectors of Σ . Let $\lambda_1, \dots, \lambda_N$ denote the diagonal terms of Λ , i.e. the eigenvalues of Σ . Let $C_* = C(C'\Sigma^{-1}C)^{-1}\gamma$. It is the linear combination of constraints where the coefficient of each constraint is its shadow price. Then Equation (1.3) can be rewritten as $w = \Sigma^{-1}C_* = U\Lambda^{-1}U'C_*$, or as:

$$w = \sum_{i=1}^N \frac{C'_* u_i}{\lambda_i} u_i. \quad (1.4)$$

The constrained minimum variance portfolio spreads its weight across the eigenvectors of Σ . The weight on eigenvector u_i is inversely proportional to its eigenvalue λ_i . λ_i is the variance of returns on the portfolio with weights u_i . It measures the riskiness of u_i . If an eigenvector is less risky, it receives more weight; riskier, less weight. This is the mathematical translation of the economic idea of diversification. Spreading weights across eigenvectors is like putting all the eggs in different baskets.

In practice, Σ is not known, so we can be tempted to replace it with the sample covariance matrix $\tilde{\Sigma}$. Decompose it into $\tilde{\Sigma} = \tilde{U}\tilde{\Lambda}\tilde{U}'$, where \tilde{U} is the rotation matrix whose columns $\tilde{u}_1, \dots, \tilde{u}_N$ are the eigenvectors of $\tilde{\Sigma}$, and $\tilde{\Lambda}$ the diagonal matrix whose diagonal terms $\tilde{\lambda}_1, \dots, \tilde{\lambda}_N$ are the eigenvalues of $\tilde{\Sigma}$. Portfolio selection with $\tilde{\Sigma}$ yields weights $\tilde{w} = \sum_{i=1}^N (\tilde{C}'_* \tilde{u}_i / \tilde{\lambda}_i) \tilde{u}_i$, where $\tilde{C}_* = C(C'\tilde{\Sigma}^{-1}C)^{-1}\gamma$.

The true riskiness of eigenvector \tilde{u}_i is $\tilde{u}_i' \Sigma \tilde{u}_i$, estimated by $\tilde{u}_i' \tilde{\Sigma} \tilde{u}_i = \tilde{\lambda}_i$. If $\tilde{\lambda}_i$ is close to zero but $\tilde{u}_i' \Sigma \tilde{u}_i$ is not, it is a catastrophe. Since weight is in $1/\tilde{\lambda}_i$, if $\tilde{\lambda}_i$ is near zero by mistake, nearly infinite weight falls on an eigenvector that is not truly riskless. It is like putting all the eggs in the same basket, and discovering that it is not safe when all the eggs get broken. A covariance matrix estimator for portfolio selection must refrain from having eigenvalues near zero, unless there is convincing evidence that it is no mistake. This is the same as saying that the covariance matrix must not be singular or even near-singular, an idea already known to Michaud (1989).

Next, we show that some eigenvalues of the sample covariance matrix are systematically

too close to zero by mistake, when N is not negligible with respect to T . The sample covariance matrix is typically singular or near-singular in practical applications. This is what makes it ill-suited to portfolio selection.

1.2.5 Sample Covariance Matrix Eigenvalues

We are trying to show that the smallest eigenvalues of the sample covariance matrix are biased towards zero. Since they are constrained to be nonnegative, we need to show that they are biased downwards. The full picture is that the smallest eigenvalues are biased downwards and the largest ones upwards. This statement is equivalent to saying that sample eigenvalues are too dispersed.

Theorem 2 *Sample eigenvalues have approximately the same average as true ones, in the sense that $E[(1/N) \sum_{i=1}^N \tilde{\lambda}_i] = (1/N) \sum_{i=1}^N \lambda_i$ and $\text{Var}[(1/N) \sum_{i=1}^N \tilde{\lambda}_i] \rightarrow 0$.*

Yin (1986) proves a more general version of this result, but under stronger assumptions. Recall from above that $m = \Sigma \circ I = (1/N) \sum_{i=1}^N \lambda_i$.

Theorem 3 *Sample eigenvalues are more dispersed than true ones, in the sense that:*

$$E \left[\frac{1}{N} \sum_{i=1}^N (\tilde{\lambda}_i - m)^2 \right] = \frac{1}{N} \sum_{i=1}^N (\lambda_i - m)^2 + E \left[\|\tilde{\Sigma} - \Sigma\|^2 \right] \quad (1.5)$$

Yin (1986) proves a related result under stronger assumptions.²

$\tilde{\Sigma}$ uses all of its error to feed an increase in the dispersion of its eigenvalues. It is as if $\tilde{\Sigma}$ wanted to have the most dispersed eigenvalues, and used all that differentiates it from Σ to beat Σ at this game. Theorem 3 implies that the smallest eigenvalues of $\tilde{\Sigma}$ are biased downwards (towards zero), and the largest ones upwards. Ironically, it is due to the fact that sample covariance matrix *entries* are *unbiased*, as is apparent from the proof of Theorem 3.

A property of eigenvalues helps understand the mechanism at work.

Theorem 4 *The eigenvalues are the most dispersed diagonal elements that can be obtained by rotation.*

²He proves that $(1/N) \sum_{i=1}^N (\tilde{\lambda}_i - m)^2 - \{(1/N) \sum_{i=1}^N (\lambda_i - m)^2 + (N/T) m^2\} \rightarrow 0$ in probability. His result follows from Theorems 1 and 3.

Since $\tilde{\Sigma}$ is unbiased and U is nonstochastic, $U'\tilde{\Sigma}U$ is an unbiased estimator of $U'\Sigma U$. The diagonal elements of $U'\tilde{\Sigma}U$ are approximately as dispersed as the ones of $U'\Sigma U$. For convenience, let us speak as if they were exactly as dispersed. By contrast, $\tilde{U}'\tilde{\Sigma}\tilde{U}$ is not an unbiased estimator of $\tilde{U}'\Sigma\tilde{U}$. This is because the errors of \tilde{U} and $\tilde{\Sigma}$ strongly interact. By Theorem 4, the diagonal elements of $\tilde{U}'\tilde{\Sigma}\tilde{U}$ are more dispersed than those of $U'\tilde{\Sigma}U$ and $U'\Sigma U$. This is why sample eigenvalues are more dispersed than true ones.

Evidence against the sample covariance matrix is even more damning than Theorem 3 suggests, because $\tilde{\lambda}_i = \tilde{u}_i'\tilde{\Sigma}\tilde{u}_i$ should not be compared to $\lambda_i = u_i'\Sigma u_i$, but to $\tilde{u}_i'\Sigma\tilde{u}_i$. We should compare estimated vs. true riskiness of eigenvector \tilde{u}_i . In portfolio selection, we entrust our money to \tilde{u}_i based on $\tilde{u}_i'\tilde{\Sigma}\tilde{u}_i$, and we end up bearing the risk $\tilde{u}_i'\Sigma\tilde{u}_i$. By Theorem 4 again, the diagonal elements of $\tilde{U}'\tilde{\Sigma}\tilde{U}$ are even less dispersed than those of $U'\Sigma U$. Not only are sample eigenvalues more dispersed than true ones, but they should be less dispersed! Intuitively: statisticians should shy away from taking a strong stance on extremely small and large eigenvalues, because they know that they don't know everything. The sample covariance matrix is guilty of taking an unjustifiably strong stance.

How important is this effect in practice? When variables outnumber observations, it is infinitely important. Since $\tilde{\Sigma} = (1/T)XX'$ and the dimension of X is $N \times T$, the rank of $\tilde{\Sigma}$ is the minimum of N and T . When $N > T$, the rank of $\tilde{\Sigma}$ is less than its dimension N . $\tilde{\Sigma}$ is rank-deficient. This means that it is singular and that some of its eigenvalues are equal to zero. It cannot be inverted and used for portfolio selection.

By continuity, we expect the sample covariance matrix to become near-singular as the ratio N/T gets close to one. In order to see how sample covariance matrix eigenvalues change in the ratio N/T , we look more closely at a particular case. It is our experience that what follows is representative of the general case.

1.2.6 Particular Case: the Identity Matrix

To illustrate how dangerous the sample covariance matrix is for portfolio selection, we analyze in more detail the particular case $\Sigma = I$. Assuming that the ratio N/T converges to a finite positive limit c called the concentration, Marčenko and Pastur (1967) derive the

limit of the distribution of sample eigenvalues.

A popular way to graph eigenvalues is to sort them in descending order, and plot the eigenvalues as a function of their rank. We follow this convention, with one adjustment due to the fact that the number of eigenvalues goes to infinity. We plot the eigenvalues as a function of their relative rank, defined as the rank divided by the total number of eigenvalues. As N goes to infinity, the relative rank remains between zero (largest eigenvalues) and one (smallest).

By assumption, $\Sigma = I$, therefore true eigenvalues are all equal to one. Their graph is a horizontal line at one. Figure F-1 plots sample eigenvalues for various concentrations, as given by Marčenko and Pastur's asymptotic approximation. If concentration was zero, sample eigenvalues would also plot as a horizontal line at one. However, for positive concentrations, even small ones, the smallest eigenvalues are substantially biased towards zero. Bias becomes more severe as concentration increases to one. When $c > 1$, the smallest eigenvalues are *equal* to zero.

Figure F-1 speaks against using the sample covariance matrix for portfolio selection unless N is negligible with respect to T , which is rarely the case in practice. From the above discussion, it is because the sample covariance matrix uses the accumulation of errors off the diagonal to bias the smallest eigenvalues downwards and the largest ones upwards. This is a widespread phenomenon. For example, it is well-known that the smallest estimated betas are biased downwards and the largest ones upwards. It can even be said that this phenomenon plays an important role in the popularity of alternatives to the maximum likelihood such as Bayesian statistics and decision theory. It is particularly pronounced here because the excess dispersion of sample eigenvalues is in N/T , instead of e.g. $1/T$ for betas. Also, it is particularly damaging, because the downwards bias of the smallest eigenvalues, when it draws them close to zero, has infinitely destructive consequences on portfolio selection.

The bias of the eigenvalues of the sample covariance matrix is intimately related to the unbiasedness of its entries. To put it bluntly, either the eigenvalues or the entries must be biased: we cannot have it both ways. Equation (1.4) makes it clear that portfolio selection calls for minimally biased eigenvalues, even if the price to pay is to bias the entries. This

is the topic of Section 1.3.

1.2.7 Potential Applications to Tests for the Number of Factors in the APT

Some of the plots in Figure F-1 bear a striking resemblance to plots of the eigenvalues of the sample covariance matrix of stock returns in tests for the number of factors in the APT. There, the emphasis is not on the smallest eigenvalues, but on the largest ones: are they large enough to support the APT? As can be seen from Figure F-1, the largest sample eigenvalues are severely biased upwards, therefore inference must be drawn cautiously. This is the point made by Brown (1989), based on Monte-Carlo simulations. The review by Connor and Korajczyk (1992) makes it clear that this is a pervasive problem in the literature.

Marčenko and Pastur (1967) solve much more than the special case $\Sigma = I$. They derive a general equation that yields the distribution of sample eigenvalues as a function of the distribution of true eigenvalues and the concentration. An original approach to APT tests would be to use this equation in reverse to back up true eigenvalues from sample eigenvalues. This is an appealing direction for future research, but there is one obstacle. It is an ill-posed problem.

Infinitesimal errors on the estimation of sample eigenvalues are amplified into large errors on true eigenvalues as we go through the equation in reverse. For example, Black and Scholes (1973) obtain a partial differential equation that determines the value $V(S, t)$ of a European option as a function of the stock price S and time t . They know $V(\cdot, t_2)$ at expiration date t_2 , and want $V(\cdot, t_1)$ today at $t_1 < t_2$. This is a well-posed problem. Reverse the direction of time and it becomes an ill-posed problem. It would not be possible to deduce $V(\cdot, t_2)$ from $V(\cdot, t_1)$ for $t_2 > t_1$. More precisely, a lot of very different solutions $V(\cdot, t_2)$ correspond to almost exactly the same initial conditions $V(\cdot, t_1)$. Fortunately for option pricing, time flows in the right direction.

The distribution of sample eigenvalues is a smoothed-out version of the distribution of true eigenvalues. It is a general fact that “un-smoothing” is an ill-posed problem. Figuratively, this is because the resolution of the picture is diminished by the action of

smoothing. In our case, it is the error of sample eigenvectors that smoothes out true eigenvalues into sample eigenvalues. For option pricing, it is the uncertainty about the terminal value of the stock price that makes today's option value $V(\cdot, t_1)$ smoother than the terminal payoff $V(\cdot, t_2)$.

Ill-posedness makes it hard to obtain reliable estimators of true eigenvalues. Getting confidence intervals is probably even harder. Not surprisingly, the degree of ill-posedness increases in the ratio N/T . We interpret it as: we cannot get something for nothing. We firmly believe that ill-posedness is not an artifact of the Marčenko and Pastur equation, but a deep feature of the problem itself.

However, the degree of ill-posedness is not uniform. The problem is better posed around isolated eigenvalues. In practice, we expect the largest eigenvalues to be quite isolated. This may be what makes it possible to recover them. Some more details are in Appendix A. For a different and innovative approach, see Adamek (1994).

1.3 Improved Covariance Matrix Estimation

We derive an estimator that improves over the sample covariance matrix when the number of variables N is not negligible with respect to the number of observations T . Generalizations are described.

1.3.1 Linear Shrinkage of Sample Eigenvalues

As we saw in Section 1.2, the problem with the sample covariance matrix is that its eigenvalues can be too dispersed. The line of attack is suggested by established methods in multivariate statistics. Muirhead (1987) reviews decision-theoretic alternatives to the sample covariance matrix and concludes that they “have a tendency to move the sample eigenvalues together in an intuitively appealing way.” Shrinking sample eigenvalues together is attractive for portfolio selection because it reduces singularity by pulling the smallest eigenvalues away from zero. We follow this approach.

To simplify matters, we focus on linear shrinkage. That is, we consider improved

eigenvalues estimators of the form $\hat{\lambda}_i = \alpha + \beta \tilde{\lambda}_i$, $i = 1, \dots, N$, where α and β are scalars.³ This is equivalent to replacing $\tilde{\Lambda}$ with $\hat{\Lambda} = \alpha I + \beta \tilde{\Lambda}$. Following the decision-theoretic literature, we keep the same eigenvectors as the sample covariance matrix. The improved estimator is: $\hat{\Sigma} = \tilde{U} \hat{\Lambda} \tilde{U}' = \tilde{U}(\alpha I + \beta \tilde{\Lambda}) \tilde{U}' = \alpha I + \beta \tilde{\Sigma}$.

The central question is to find the coefficients α and β . If we were only trying to avoid singularity, the choice of α and β would be ad-hoc. Instead, we ought to be minimizing some criterion. A natural candidate is the mean squared error:

$$\begin{aligned} \min_{\alpha, \beta} E \left[\left\| \hat{\Sigma} - \Sigma \right\|^2 \right] \\ \text{s.t. } \hat{\Sigma} = \alpha I + \beta \tilde{\Sigma}. \end{aligned} \quad (1.6)$$

Is it compatible with the need to avoid singularity? $E[\|\hat{\Sigma} - \Sigma\|^2] = E[\|\tilde{U}' \hat{\Sigma} \tilde{U} - \tilde{U}' \Sigma \tilde{U}\|^2] = E[(1/N) \sum_{i=1}^N (\tilde{u}_i' \hat{\Sigma} \tilde{u}_i - \tilde{u}_i' \Sigma \tilde{u}_i)^2] + \text{constant}$, where the constant does not depend on α and β . Therefore choosing α and β to minimize mean squared error is the same as choosing them to minimize the distance between the estimated riskiness $\tilde{u}_i' \hat{\Sigma} \tilde{u}_i$ of eigenvector \tilde{u}_i and its true riskiness $\tilde{u}_i' \Sigma \tilde{u}_i$, on average across $i = 1, \dots, N$. For portfolio selection, this is a very good criterion, since so much rides on estimating the riskiness of each eigenvector well. The mean squared error criterion is in alignment with the objectives of portfolio selection. Even more alignment could conceivably be achieved, for example by letting the criterion depend on the matrix of portfolio selection constraints C (cf. Equation (1.2)), but this is left to future research.⁴

$\tilde{u}_1' \Sigma \tilde{u}_1, \dots, \tilde{u}_N' \Sigma \tilde{u}_N$ are even less dispersed than true eigenvalues, so we anticipate that our estimator's eigenvalues will be less dispersed than true ones. This should keep the smallest eigenvalues of $\hat{\Sigma}$ safely away from zero.

1.3.2 Optimal Linear Shrinkage

If we could observe the true covariance matrix Σ , we could easily solve Equation (1.6).

³One advantage of linear shrinkage is that it preserves the ordering of the eigenvalues (if $\beta \geq 0$), an intuitively appealing property whose theoretical importance is proven by Sheena and Takemura (1992).

⁴I thank Fischer Black for this suggestion.

Theorem 5 Let $m = \Sigma \circ I$. Let $r_1^2 = \|\Sigma - mI\|^2$, $r_2^2 = E[\|\tilde{\Sigma} - \Sigma\|^2]$ and $d^2 = E[\|\tilde{\Sigma} - mI\|^2]$. The solution $\hat{\Sigma}$ to Equation (1.6) is:

$$\hat{\Sigma} = \frac{r_2^2}{d^2} mI + \frac{r_1^2}{d^2} \tilde{\Sigma}. \quad (1.7)$$

Its mean squared error is $E[\|\hat{\Sigma} - \Sigma\|^2] = r_1^2 r_2^2 / d^2 < \min(r_1^2, r_2^2)$.

By Theorem 3, $r_1^2 + r_2^2 = d^2$, so $\hat{\Sigma}$ is a weighted average of mI and $\tilde{\Sigma}$. The weight placed on mI increases with the error of $\tilde{\Sigma}$ and decreases with the error of mI . For the weight on $\tilde{\Sigma}$, it is the opposite. The dispersion of the eigenvalues of $\hat{\Sigma}$ is $E[\|\hat{\Sigma} - mI\|^2] = r_1^4 / d^2 < r_1^2$: the eigenvalues of $\hat{\Sigma}$ are even less dispersed than Σ 's. This effect becomes more pronounced as the error of $\tilde{\Sigma}$ increases, i.e. as the ratio N/T increases. $\hat{\Sigma}$ is the projection of Σ onto the line between mI and $\tilde{\Sigma}$. Figure F-2 shows this geometrical interpretation.

Unfortunately, $\hat{\Sigma}$ is not an estimator because it depends on the unobservable matrix Σ . As we saw, in general it is impossible to estimate Σ consistently. However, we do not need all the entries of Σ : the four parameters m , r_1^2 , r_2^2 and d^2 suffice. The key insight of this chapter is that, as T goes to infinity, even if N goes to infinity too, it is possible to estimate these four parameters consistently.

First, Theorem 2 reveals that m can be estimated simply by $\hat{m} = (1/N) \sum_{i=1}^N \tilde{\lambda}_i$: the average of sample eigenvalues is a consistent estimator of the average of true eigenvalues. Second, a natural estimator of $d^2 = E[\|\tilde{\Sigma} - mI\|^2]$ is $\hat{d}^2 = \|\tilde{\Sigma} - \hat{m}I\|^2$.

Theorem 6 $\hat{d} - d \xrightarrow{P} 0$, where \xrightarrow{P} denotes convergence in probability as T goes to infinity.

Third, let the $N \times 1$ vector x_t denote the t^{th} column of the observations matrix X for $t = 1, \dots, T$. $\tilde{\Sigma} = (1/T)XX'$ can be rewritten as $\tilde{\Sigma} = (1/T) \sum_{t=1}^T x_t x_t'$. $\tilde{\Sigma}$ is the average of the matrices $x_t x_t'$ ($t = 1, \dots, T$). Since the matrices $x_t x_t'$ are iid across $t = 1, \dots, T$, we can estimate the error $d^2 = E[\|\tilde{\Sigma} - \Sigma\|^2]$ of their average by seeing how far each one of them deviates from the average.

Theorem 7 Define $\hat{r}_2^2 = (1/T^2) \sum_{t=1}^T \|x_t x_t' - \tilde{\Sigma}\|^2$. Then $\hat{r}_2^2 - r_2^2 \xrightarrow{P} 0$.

Finally, Theorem 3 can be rewritten as $r_1^2 = d^2 - r_2^2$.

Theorem 8 Define $\hat{r}_1^2 = \hat{d}^2 - \hat{r}_2^2$. Then $\hat{r}_1^2 - r_1^2 \xrightarrow{P} 0$.

If, for a given realization, $\hat{d}^2 < \hat{r}_2^2$, then we recommend correcting \hat{d}^2 and/or \hat{r}_2^2 so that they are equal. It can be shown that this does not affect the validity of the theorems.

Please note that Theorems 6, 7 and 8 are non-trivial since, in spite of the division by N in the definition of the norm $\|\cdot\|$, the scalars d , r_1 , and r_2 do not converge to zero (except in special cases), as is apparent from the proofs.

Plugging consistent estimators in place of the unobservable parameters in Equation (1.7) yields a consistent estimator of $\hat{\Sigma}$ with the same asymptotic properties. This is the main result of the chapter.

Theorem 9 *The improved estimator*

$$\hat{\hat{\Sigma}} = \frac{\hat{r}_2^2}{\hat{d}^2} \hat{m}I + \frac{\hat{r}_1^2}{\hat{d}^2} \hat{\Sigma} \quad (1.8)$$

estimates the solution $\hat{\Sigma}$ to Equation (1.6) consistently, i.e. $\|\hat{\hat{\Sigma}} - \hat{\Sigma}\|^2 \xrightarrow{P} 0$. Both $\hat{\hat{\Sigma}}$ and $\hat{\Sigma}$ have the same asymptotic mean squared error, i.e. $E[\|\hat{\hat{\Sigma}} - \Sigma\|^2] - E[\|\hat{\Sigma} - \Sigma\|^2] \rightarrow 0$, and $\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2$ estimates it consistently, i.e. $(\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2) - (r_1^2 r_2^2 / d^2) \xrightarrow{P} 0$.

$\hat{\hat{\Sigma}}$ is an improved estimator of the covariance matrix. It is a consistent estimator of the linear combination of the sample covariance matrix with the identity matrix that minimizes mean squared error. It is easy to verify that $\hat{\hat{\Sigma}}$ is invariant by rotation, i.e. premultiplying the observations X by a rotation matrix V ($V'V = VV' = I$) changes $\hat{\hat{\Sigma}}$ into $V'\hat{\hat{\Sigma}}V$.

By Theorem 1, the weight on $\hat{m}I$ increases in N/T . If N remains bounded, asymptotically all the weight is on the sample covariance matrix $\hat{\Sigma}$.

The advantage of our framework over finite sample statistics is that we do not have to take into account the error of estimators of the unobservable parameters m , r_1^2 , r_2^2 and d^2 . The advantage over standard asymptotics is that we encompass realistic situations where the sample covariance matrix is not optimal.

1.3.3 Generalization

$\hat{\Sigma}$ is a weighted average of $\hat{m}I$ and $\tilde{\Sigma}$. $\hat{m}I$ can be thought of as an estimator of the covariance matrix. It has asymptotically minimum mean squared error among a certain class of estimators: scalar multiples of the identity matrix. This class imposes a lot of structure on the covariance matrix: no covariances, and all variances are the same. There is only one free parameter, as opposed to $N(N + 1)/2$ for $\tilde{\Sigma}$. This parsimonious structure makes $\hat{m}I$ heavily biased, but at least it prevents it from being singular, a problem that hurts the unstructured, unbiased estimator $\tilde{\Sigma}$.

Other structures can be imposed on the covariance matrix. Frost and Savarino (1986) impose that all stock returns have the same variance and all pairs of stock returns have the same covariance. They have two free parameters. We can also impose that the covariance matrix is diagonal (N parameters), or that all correlation coefficients are equal ($N + 1$ parameters).

We call such estimators: “structured.” Other structured estimators of interest in Finance are the index models. For example, Sharpe’s (1963) single index model assumes that the idiosyncratic risks of different stocks are uncorrelated. The idiosyncratic risk is the fraction of the risk that is not systematic risk. Systematic risk is the fraction of the risk that can be explained as covariance with an index, usually a broad-based market index. In general, if there are K indices, then we need to estimate the covariance matrix of the indices ($K(K + 1)/2$ parameters), the covariance of each stock with each index (KN parameters) and each stock’s idiosyncratic risk (N parameters), for a total of $(K + 1)(N + K/2)$ free parameters. When $K \ll N$, this is still much fewer parameters than the sample covariance matrix.

Structured estimators are popular for portfolio selection. They are carefully designed to avoid the singularity problem of the sample covariance matrix. Their main selling point is that they do not place infinite weights on risky eigenvectors by mistake.

However, the way that they obtain this desirable feature is ad-hoc. They impose arbitrary structure that they know is wrong, then disregard any evidence that goes against it. They throw away all sample information that does not fit in their arbitrarily specified structure.

It would be better to recycle the information that they ignore, in an optimal way. We recommend taking a well-chosen weighted average of a structured estimator and the sample covariance matrix.

Let $\bar{\Sigma}$ denote any given structured estimator of interest to the statistician. Consider the problem:

$$\begin{aligned} \min_{\omega} E \left[\|\hat{\Sigma} - \Sigma\|^2 \right] \\ \text{s.t. } \hat{\Sigma} = \omega \bar{\Sigma} + (1 - \omega) \tilde{\Sigma}. \end{aligned} \quad (1.9)$$

$\hat{\Sigma}$ is a weighted average of two estimators, one generally singular ($\tilde{\Sigma}$), and the other one generally not ($\bar{\Sigma}$). Which one does it inherit its properties from? An elementary result from matrix algebra answers.

Proposition 1 *The smallest eigenvalue of $\hat{\Sigma} = \omega \bar{\Sigma} + (1 - \omega) \tilde{\Sigma}$ is at least as large as ω times the smallest eigenvalue of $\bar{\Sigma}$.*

$\bar{\Sigma}$ is constructed so that its smallest eigenvalues do not come near zero. Therefore $\hat{\Sigma}$ is generally not singular, unless ω is very small. If ω was very small, then it would mean that the sample covariance matrix can hardly be improved on. From what we have seen so far, this would be rather surprising when N is of the same order of magnitude as T .

Again, let us pretend for a moment that we can observe Σ . As above, let $r_1^2 = E[\|\bar{\Sigma} - \Sigma\|^2]$, $r_2^2 = E[\|\tilde{\Sigma} - \Sigma\|^2]$ and $d^2 = E[\|\tilde{\Sigma} - \bar{\Sigma}\|^2]$. In addition, let $\varphi = E[(\bar{\Sigma} - \Sigma) \circ (\tilde{\Sigma} - \Sigma)]$ measure the “covariance” between the errors of both estimators.

Theorem 10 *Then the solution to Equation (1.9) is given by:*

$$\hat{\Sigma} = \frac{r_2^2 - \varphi}{d^2} \bar{\Sigma} + \left(1 - \frac{r_2^2 - \varphi}{d^2} \right) \tilde{\Sigma}. \quad (1.10)$$

The geometric interpretation is the same as in Figure F-2, except that $\bar{\Sigma}$ replaces mI and that the triangle $(\bar{\Sigma}, \Sigma, \tilde{\Sigma})$ does not necessarily have a right angle at Σ anymore. In the particular case $\varphi = 0$, the weight on $\bar{\Sigma}$ reduces to r_2^2/d^2 , as in Equation (1.7). This simplification takes place (asymptotically) for $\bar{\Sigma} = \widehat{m}I$, but not necessarily for other structured estimators.

Again the problem is to estimate the unobservable parameters r_1^2 , r_2^2 , d^2 and φ consistently. We do not provide formal proofs of consistency, since they would have to be

rewritten for every structured estimator $\bar{\Sigma}$. We just indicate how the general logic of the argument for $\bar{\Sigma} = \hat{m}I$ can be extended to other structured estimators. In Section 1.4, we provide empirical support for these extensions.

We can take the same estimator \hat{r}_2^2 as before. The estimator of d^2 becomes $\hat{d}^2 = \|\hat{\Sigma} - \bar{\Sigma}\|^2$. The additional complication is that we need an estimator $\hat{\varphi}$ of φ . Since $d^2 = r_1^2 + r_2^2 - 2\varphi$, $\hat{\varphi}$ would let us estimate r_1^2 by $\hat{r}_1^2 = \hat{d}^2 - \hat{r}_2^2 + 2\hat{\varphi}$.

Let $\bar{\Sigma} = [\bar{\sigma}_{ij}]_{i,j=1,\dots,N}$ and $\tilde{\Sigma} = [\tilde{\sigma}_{ij}]_{i,j=1,\dots,N}$. Since $\varphi = (1/N) \sum_{i=1}^N \sum_{j=1}^N \text{Cov}[\bar{\sigma}_{ij}, \tilde{\sigma}_{ij}]$, all we need is estimators of $\varphi_{ij} = \text{Cov}[\bar{\sigma}_{ij}, \tilde{\sigma}_{ij}]$ for $i, j = 1, \dots, N$. They are usually suggested by the nature of $\bar{\Sigma}$. The idea is that, if we can estimate $\bar{\sigma}_{ij}$, then we can estimate the error on $\bar{\sigma}_{ij}$, and its covariance with the error on $\tilde{\sigma}_{ij}$. Please keep in mind that φ_{ij} vanishes in $1/T$, even though φ itself may be of order N/T . Therefore, in the more complicated cases, the delta method can be used to estimate φ_{ij} consistently. Given the estimators $\hat{\varphi}_{ij}$ for $i, j = 1, \dots, N$, we form $\hat{\varphi} = (1/N) \sum_{i=1}^N \sum_{j=1}^N \hat{\varphi}_{ij}$.

Appendix B gives the formula of $\hat{\varphi}_{ij}$ ($i, j = 1, \dots, N$) for various structured estimators.

1.3.4 Comparison with Previous Work in Multivariate Statistics

This approach has an obvious Bayesian interpretation. Bayesian statistics combine sample information with other sources of information. The other sources are summarized in a “prior” distribution of the unknown parameter. In our case, the prior distribution puts all its mass on a sphere centered on $\bar{\Sigma}$ with radius \hat{r}_1 . Then sample information reveals that the true parameter also lies on another sphere, with center $\tilde{\Sigma}$ and radius \hat{r}_2 . Combining prior and sample yields a posterior distribution. In our case, the true covariance matrix must lie on the intersection of the two spheres. This intersection is a circle. At the center of this circle stands the improved estimator $\hat{\hat{\Sigma}}$. This interpretation is shown in Figure F-3.

Fundamental Bayesian questions are: Where does the prior come from? How confident are we in the prior? In finite sample, it is very hard to answer these questions satisfactorily. If the statistician chooses the prior without looking at data at all, it might be very inaccurate. Empirical Bayesians do look at data, but then they pretend that they did not, and ignore dependence between prior and sample. In some cases, dependence can safely be neglected,

but how do we know that?

By contrast, in our asymptotic framework, we can build the prior around any structured estimator already used in practice. Furthermore, the degree of confidence in the prior can be estimated consistently. In particular, we estimate the parameter φ that captures dependence between prior and sample. We find out for any given prior whether φ can be neglected, and if it cannot be, we account for it in Equation (1.10).

In the established nomenclature, our work is not pure Bayesian because we estimate the prior from the sample. It is not empirical Bayesian either because it takes into account the dependence between the estimated prior and the sample. It is decision theory.

For the covariance matrix, previous literature on decision theory (and on pure and empirical Bayesian statistics too) has been only in finite sample. The reason is that, under standard asymptotics, the sample covariance matrix is consistent, so there is no need to seek alternatives. Decision theory in finite sample is not very tractable. Also, it relies on the Wishart distribution, which has two limitations: random variables must be normally distributed, and if variables outnumber observations then the Wishart density does not exist (because $\tilde{\Sigma}$ is rank-deficient). For portfolio selection, both limitations are serious.

One of our contributions is to realize that these are not limitations of decision theory itself, but of finite sample. In stock market finance, we are fortunate enough to have large numbers of observations, which make asymptotic approximations realistic, and large numbers of variables, which open the door to improvements over the sample covariance matrix. This is the ideal situation to free decision theory from finite sample drawbacks. All that is needed is to relax the standard asymptotic assumption that keeps the number of variables bounded. $\hat{\Sigma}$ is the first estimator of the covariance matrix based on *asymptotic* decision theory.

Stein (1975) suggests that invariance by rotation is an important property for covariance matrix estimators. Intuitively, it means that the statistician lets the data speak without putting a spin on what they say. This excludes all of the structured estimators cited above except $\widehat{m}I$. The existing literature does not contain any estimator invariant by rotation and

theoretically motivated when $N > T$. Perhaps more importantly, it contains no estimator that is invariant by rotation and is known not to be singular or near-singular when $N > T$. This has lead some to believe that the inverse of the covariance matrix could not be estimated at all when $N > T$.

Now it can be. The estimator $\hat{\hat{\Sigma}}$ of Section 1.3.2 is invariant by rotation. It has a sound theoretical motivation when $N > T$. As a matter of fact, it does not even matter whether $N > T$, which is satisfying because we should expect some continuity between $N = 999$, $T = 1000$ and $N = 1000$, $T = 999$. The eigenvalues of $\hat{\hat{\Sigma}}$ are asymptotically even less dispersed than Σ 's, which prevents $\hat{\hat{\Sigma}}$ from being near-singular or singular. The dispersion of the eigenvalues of $\hat{\hat{\Sigma}}$ actually *decreases* in the ratio N/T . Therefore $\hat{\hat{\Sigma}}^{-1}$ is the first estimator of inverse of the covariance matrix that is invariant by rotation and can be used when variables outnumber observations.

1.4 Application to Portfolio Selection

The goal of this section is to find out how the asymptotic results of Section 1.3 carry through to large but finite sample. We first compare $\hat{\hat{\Sigma}}$ to other estimators in terms of mean squared error in Monte-Carlo simulations. Then we apply $\hat{\hat{\Sigma}}$ to historical stock returns data.

1.4.1 Monte-Carlo Simulations

Our purpose is to compare the mean squared errors of various estimators across a range of situations. We focus on estimators that are invariant by rotation, therefore we use Equation (9) for $\hat{\hat{\Sigma}}$.

The benchmark is the mean squared error of the covariance matrix. We report the Percentage Relative Improvement in Average Loss of $\hat{\hat{\Sigma}}$, defined as: $\text{PRIAL}(\hat{\hat{\Sigma}}) = (E[\|\tilde{\Sigma} - \Sigma\|^2] - E[\|\hat{\hat{\Sigma}} - \Sigma\|^2]) / E[\|\tilde{\Sigma} - \Sigma\|^2] \times 100$. If the PRIAL is positive (negative), then $\hat{\hat{\Sigma}}$ performs better (worse) than $\tilde{\Sigma}$. The PRIAL of the sample covariance matrix is zero by definition. The PRIAL cannot exceed 100%. We compare the PRIAL of $\hat{\hat{\Sigma}}$ to the PRIAL of other estimators from finite sample decision theory.

Haff (1980) introduces an estimator with an empirical Bayesian interpretation. Like $\hat{\hat{\Sigma}}$, it

is a linear combination of the sample covariance matrix and the identity. The difference lies in the coefficients of the combination. Haff's coefficients do not depend on the observations X , only on N and T . If the criterion is the mean squared error, Haff's approach suggests:

$$\hat{\Sigma}_{EB} = \frac{NT - 2T - 2}{NT^2} \widehat{m}_{EB} I + \frac{T}{T+1} \tilde{\Sigma} \quad (1.11)$$

with $\widehat{m}_{EB} = [\det(\tilde{\Sigma})]^{1/N}$. When $N > T$ we take $\widehat{m}_{EB} = \widehat{m}$ because the regular formula would yield zero. The initials EB stand for empirical Bayesian.

Stein (1975) proposes an estimator that keeps the eigenvectors of the sample covariance matrix and replaces its eigenvalues $\tilde{\lambda}_1, \dots, \tilde{\lambda}_N$ by:

$$T\tilde{\lambda}_i / \left(T - N + 1 + 2\tilde{\lambda}_i \sum_{\substack{j=1 \\ j \neq i}}^N \frac{1}{\tilde{\lambda}_i - \tilde{\lambda}_j} \right) \quad i = 1, \dots, N. \quad (1.12)$$

These corrected eigenvalues need neither be positive nor in the same order as sample eigenvalues. To prevent this from happening, an ad-hoc procedure called isotonic regression is applied before recombining corrected eigenvalues with sample eigenvectors.⁵ Haff (1982) independently obtains a closely related estimator. In any given simulation, we call $\hat{\Sigma}_{SH}$ the better performing estimator of the two. The other one is not reported. The initials SH stand for Stein and Haff.⁶

Stein (1982) and Dey and Srinivasan (1985) both derive the same estimator. Under a certain loss function, it is minimax, which means that no other estimator has lower worst-case error. The minimax criterion is sometimes criticized as overly pessimistic, since it looks at the worst case only. This estimator preserves sample eigenvectors and replaces sample eigenvalues by:

$$\frac{T}{T + N + 1 - 2i} \tilde{\lambda}_i, \quad (1.13)$$

where sample eigenvalues $\tilde{\lambda}_1, \dots, \tilde{\lambda}_N$ are sorted in descending order. We call this estimator

⁵Intuitively, isotonic regression restores the ordering by assigning the same value to a subsequence of corrected eigenvalues that would violate it.

⁶When $N > T$ some of the terms $\tilde{\lambda}_i - \tilde{\lambda}_j$ in formula (1.12) result in a division by zero. We just ignore them. Nonetheless, when N is too large compared to T , the isotonic regression does not converge. In this case $\hat{\Sigma}_{SH}$ does not exist.

$\hat{\Sigma}_{MX}$. The initials MX stand for minimax.

We simulate normally distributed random variables. The true covariance matrix Σ can be taken diagonal without loss of generality. We draw its eigenvalues according to a log-normal distribution. We set their average equal to one without loss of generality. We let their dispersion r_1^2 vary around the central value $1/2$. We let the ratio N/T vary around the central value $1/2$. Finally, we let the product NT vary around the central value 800. We study the influences of r_1^2 , N/T and NT separately. When one parameter moves, the other two remain fixed at their central values.

The asymptotic PRIAL of $\hat{\hat{\Sigma}}$ implied by Theorems 1 and 9 is $(N/T)/[(N/T) + r_1^2] \times 100$. The PRIAL increases in N/T and decreases in r_1^2 . This is intuitive because N/T is the error on $\hat{\Sigma}$ and r_1^2 is the error on $\hat{m}I$.

When all three parameters are fixed at their central values, we get the results in Table 1.1. “Risk” means the average mean squared error over 1,000 simulations. For the central values

Estimator	$\tilde{\Sigma}$	$\hat{\hat{\Sigma}}$	$\hat{\Sigma}_{EB}$	$\hat{\Sigma}_{SH}$	$\hat{\Sigma}_{MX}$
Risk	0.5372	0.2723	0.5120	0.3076	0.3222
Standard Error on Risk	(0.0033)	(0.0013)	(0.0031)	(0.0014)	(0.0014)
PRIAL	0.00%	49.31%	4.69%	42.74%	40.02%

Table 1.1: Result of 1,000 Monte-Carlo Simulations for Central Parameter Values.

of the parameters, the asymptotic PRIAL of $\hat{\hat{\Sigma}}$ is 50%. Table 1.1 shows that asymptotic behavior is practically attained for $N = 20$ and $T = 40$. $\hat{\hat{\Sigma}}$ improves substantially over $\tilde{\Sigma}$ and $\hat{\Sigma}_{EB}$ and moderately over $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$. This may be due to the fact that $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$ were originally derived under another loss function than the mean squared error.

When we increase N/T from zero to infinity, the asymptotic PRIAL of $\hat{\hat{\Sigma}}$ increases from 0% to 100% with an “S” shape. Figure F-4 confirms this.⁷ $\hat{\hat{\Sigma}}$ always has lower mean squared error than $\tilde{\Sigma}$ and $\hat{\Sigma}_{EB}$. It usually has slightly lower mean squared error than $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$. $\hat{\Sigma}_{SH}$ is not defined for high values of N/T . $\hat{\Sigma}_{MX}$ performs slightly better than $\hat{\hat{\Sigma}}$ for the highest values of N/T . This may be due to the fact that $\hat{\hat{\Sigma}}$ does not attain its asymptotic

⁷Corresponding tables of results are available from the author upon request. Standard errors on our estimators of the mean squared error have the same order of magnitude as in Table 1.1.

performance for values of T below 10.

When we increase r_1^2 from zero to infinity, the asymptotic PRIAL of $\hat{\hat{\Sigma}}$ decreases from 100% to 0% with a reverse “S” shape. Figure F-5 confirms this. $\hat{\hat{\Sigma}}$ has lower mean squared error than $\tilde{\Sigma}$ always, and than $\hat{\Sigma}_{EB}$ almost always. $\hat{\hat{\Sigma}}$ always has lower mean squared error than $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$. When r_1^2 gets too large, $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$ perform worse than the sample covariance matrix. The reason is that $\hat{m}I$ is very erroneous, and they shrink sample eigenvalues together too much. It is very reassuring that, in a case where its leading competitors perform much worse than $\tilde{\Sigma}$, $\hat{\hat{\Sigma}}$ performs at least as well as $\tilde{\Sigma}$.

When we increase NT from zero to infinity, we should see the PRIAL of $\hat{\hat{\Sigma}}$ converge to its asymptotic value of 1/2. Figure F-6 confirms this. $\hat{\hat{\Sigma}}$ always has lower mean squared error than $\tilde{\Sigma}$ and $\hat{\Sigma}_{EB}$. It has moderately lower mean squared error than $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$, except when T is below 20. When T is below 20, $\hat{\hat{\Sigma}}$ performs slightly worse than $\hat{\Sigma}_{SH}$ and moderately worse than $\hat{\Sigma}_{MX}$, but still substantially better than $\tilde{\Sigma}$.

When the number of variables N is large, $\hat{\hat{\Sigma}}$ and $\tilde{\Sigma}$ take much less time to compute than $\hat{\Sigma}_{EB}$, $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$ because they do not need eigenvalues and determinants. Indeed the number and nature of operations needed to compute $\hat{\hat{\Sigma}}$ are of the same order as for $\tilde{\Sigma}$. It can be an enormous advantage when the covariance matrix is very large. The only seemingly slow step is the estimation of r_2^2 , but it can be accelerated by writing:

$$\hat{r}_2^2 = \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N \left[\frac{1}{T} (X^{\wedge 2}) (X^{\wedge 2})' \right]_{ij} - \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N \left[\left(\frac{1}{T} X X' \right)^{\wedge 2} \right]_{ij}$$

where $[\cdot]_{ij}$ denotes the entry (i, j) of a matrix and the symbol \wedge denotes elementwise exponentiation, i.e. $[S^{\wedge 2}]_{ij} = ([S]_{ij})^2$ for any matrix S .

Simulations not reported here study departures from normality. These departures have little impact on the above results. In relative terms, $\tilde{\Sigma}$ and $\hat{\Sigma}_{EB}$ appear to suffer the most; then $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$; and $\hat{\hat{\Sigma}}$ appears to suffer the least.

We draw the following conclusions from these simulations. The asymptotic theory developed in Sections 1.2-1.3 approximates finite sample behavior well, as soon as T and N become of the order of 20 to 40. $\hat{\hat{\Sigma}}$ has lower mean squared error than the sample covariance matrix across the wide range of simulations studied. $\hat{\hat{\Sigma}}$ usually improves over

existing finite sample decision theory estimators, in terms of mean squared error.⁸ It sometimes performs substantially better than them. It never performs substantially worse than them.

This set of simulations indicates that the estimator $\hat{\Sigma}$ from Section 1.3.2 can be used as an all-purpose estimator of the covariance matrix.

1.4.2 Historical Data

This section takes the covariance matrix estimator $\hat{\Sigma}$ to the data. The objective is to estimate how well it would have performed over the past, had it been used for portfolio selection.

Monthly stock returns in excess of the riskfree rate and capitalizations from July 1926 to June 1993 are drawn from the Center for Research in Security Prices (CRSP) database. Let y denote any year between 1936 and 1992. Stock returns from July of year $y - 10$ to June of year y are used to estimate the covariance matrix of stock returns. Stocks with missing observations are excluded. We consider only common stocks traded on the New York Stock Exchange (NYSE) or the American Stock Exchange (AMEX).⁹ We require stocks to have a valid market capitalization in June of year y .

From these data, we extract two factors that past research has associated with stock returns. The first one is the beta with respect to a CRSP value-weighted index including dividends.¹⁰ The second factor is the logarithm in base 10 of the market capitalization in dollars of a given stock, minus the average logarithm of market capitalization across all stocks in the dataset in year y . We call this factor: “size”, for brevity. The average is subtracted because a stock with the same 50 million dollars capitalization would have been relatively large in 1936, and relatively small in 1992. Thus, a stock with “size” one (respectively minus one) is ten times larger (respectively smaller) than the average stock in the market.

⁸We acknowledge that $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$ were designed with another criterion than the mean squared error in mind. Our conclusions say nothing about performance under any other criterion. Nevertheless, the mean squared error is an important criterion. Also, there is some similarity between criteria, as suggested by the fact that $\hat{\Sigma}_{SH}$ and $\hat{\Sigma}_{MX}$ do perform well in terms of mean squared error.

⁹AMEX stocks do not appear in the CRSP database before July 1963. We do not include them before $y = 1973$.

¹⁰Before July 1963, the NYSE index; afterwards, the NYSE and AMEX index.

We compare different covariance matrices, either of the structured type ($\bar{\Sigma}$) or the asymptotic shrinkage type ($\hat{\hat{\Sigma}}$). We do not include the other estimators because they are too costly to compute or not defined at all when N is much larger than T , which is the case here.¹¹ $\bar{\Sigma}$ can be either $\widehat{m}I$ as in Section 1.3.2, or any of the four structured estimators in Appendix B. Each of these five structured estimators gives a shrinkage estimator. Therefore there are ten estimators in total.

We impose different sets of portfolio constraints. We always make weights sum up to one. In addition, we impose zero, one or two constraints chosen among the following two: the weighted average of betas has a required value; the weighted average of sizes has a required value. Therefore there are four possible sets of constraints.

Based on these data, we buy at the end of June of year y forty different kinds of minimum variance portfolios corresponding to the ten covariance matrices and the four sets of constraints. We hold them until the end of June of year $y + 1$, at which time they are rebalanced in a similar fashion, incorporating fresh data. This scheme yields a time series of monthly returns for each of the forty kinds of portfolios from July 1936 to June 1993. Since each rebalancing is based only on information that is available at the time, we are simulating realistic investment strategies. Tests based on strategies such as these ones, i.e. that do not require hindsight, are called predictive. They are easier to interpret than non-predictive tests. In addition, since we measure true buy-and-hold returns and rebalance portfolios only once a year, transactions costs are quite limited. We ignore them.

The most urgent questions concern shrinkage weight $(r_2^2 - \varphi)/d^2$: Is it between zero and one? Is it relatively stable over time? Does it make intuitive sense? Qualitatively, the answers to these three questions appears to be yes in Figure F-7. Weights are between 0.07 and 0.93 for every structured estimator and every year. Each structured estimator's weights remain within the same range of width 0.3 (approximately) throughout the 67 years. The ordering between weights remains the same over time, and makes intuitive sense. Diagonal structured estimators are given the least weight, probably because the true covariance matrix is far from being diagonal. The structured estimators that have the most free parameters are

¹¹The number of stocks N grows from 340 in 1936 to 1105 in 1992. The number of time periods T is 120 (ten years of monthly data).

given the highest weights, probably because they are the least biased. Qualitative evidence from Figure F-7 is very reassuring about the estimators of shrinkage weights, which are among our main contributions.

The most important question about the empirical properties of our method is: Does shrinkage help minimize variance? Table E.1 provides evidence that it does. The table shows the ex-post standard deviation of the ex-ante unconstrained minimum variance portfolio. For all five structured estimators, shrinkage yields portfolios with significantly lower variance. In some cases, variance diminishes a lot.

These results might be criticized as relying only on the unconstrained minimum variance portfolio. Therefore, for each structured estimator, we consider three portfolios: zero beta and size -1 ; unit beta and size -1 ; zero beta and unit size.¹² If an investor believed that returns are driven by beta and/or size, she would select some combination of these three portfolios. Then we give the benefit of hindsight to structured estimators, but only to them. That is, we choose the combination of these three portfolios with the lowest variance *based on ex-post variances and covariances*. We compare it to the ex-ante minimum variance portfolio from the corresponding shrinkage estimator. This is unfair because hindsight is such a strong advantage. It biases our results towards not finding that shrinkage helps reduce variance.

Results are in Table E.2. Again, all five shrinkage estimators (without hindsight) yield portfolios with lower variance than their corresponding structured estimators do (even with hindsight). In this sense, it can be said that our method yields portfolios with lower variance than could possibly be attained before. Table E.2 demonstrates empirically that our estimator $\hat{\Sigma}$ achieves its goal: it helps portfolio selection minimize variance.

Portfolios with lower variance than was previously possible open a new investment opportunity. From an economic perspective, it is interesting to know whether this new opportunity is attractive: Does it let investors improve their risk-return tradeoff? The risk-return tradeoff can be summarized by the Sharpe ratio: mean divided by standard deviation

¹²Remember that size one (minus one) means ten times larger (smaller) market capitalization than market average.

of portfolio returns.¹³

Figure F-8 plots the ex-post means and standard deviations of the ex-ante minimum variance portfolios constrained to have a specified beta between zero and one, and size zero. On each graph, portfolios obtained from a structured estimator are plotted as a dashed line, together with portfolios from the corresponding shrinkage estimator as a solid line. As seen above, the solid line ventures further into low-risk territory than the dashed line. However, the risk-return tradeoff does not seem to improve much. The dotted line, whose slope is the maximum Sharpe ratio of all the portfolios on the figure, is practically tangent to both the solid line of shrinkage estimator portfolios and the dashed line of structured estimator portfolios.

This is especially true when $\bar{\Sigma}$ is given by the single index model, which is the structured estimator closest to what actual investors would use. For the other $\bar{\Sigma}$ s, our interpretation is that combining a structured estimator with the sample covariance matrix goes a long way towards fixing its intrinsic flaws, if any exist.

Overall, the message is that low risk portfolios are penalized by low returns. They do not offer more attractive investment opportunities. While this may sound a little disappointing to a practitioner, it is on the contrary very satisfying for an economist. In equilibrium, there should be no easy and permanent way to attain an abnormally favorable risk-return tradeoff. It is rather remarkable that agents priced fairly the low-risk portfolios identified in this chapter... even long before they were identified! This can be interpreted as strong support for equilibrium theory of risk-return tradeoff.

Since a particular version of this theory, the Capital Asset Pricing Model (CAPM), has recently been challenged on empirical grounds, it is natural to extract from shrinkage covariance matrix estimators quantitative evidence on this subject beyond Figure F-8.

1.4.3 Testing an Implication of the CAPM

The CAPM implies, among other things, a positive relationship between returns and betas. A familiar method to test this is to run a cross-sectional regression of returns on betas:

¹³Returns are in excess of the riskfree rate.

the CAPM predicts a positive slope. As Fama (1970) clearly explains, this is equivalent to forming minimum variance portfolios with betas of one and zero respectively, and then testing whether they have different mean returns. This brings back CAPM tests to portfolio selection, where shrinkage covariance matrix estimators can be used.

Most existing tests run Ordinary Least Squares (OLS) regressions. This corresponds to using the structured estimator $\bar{\Sigma} = \widehat{m}I$ for portfolio selection. No doubt it can be replaced by an improved estimator of the covariance matrix. This corresponds to running Generalized Least Squares (GLS) regressions. Amihud, Christensen and Mendelson (1994) are among the few who run GLS. The problem is that they allow themselves to “peek into the future” to estimate the covariance matrix. Their test is not predictive. Its interpretation is not straightforward, because real-life investors cannot peek into the future. Furthermore, the ex-post returns that they report are not truly ex-post because they come from a period that has already been used to estimate the covariance matrix. This feature can bias standard errors towards zero, t-statistics away from zero, and tests of the CAPM towards finding a significantly positive slope. We avoid these problems by running a predictive test.

Another difficulty is estimating betas. Since beta estimates contain error, the largest ones are biased upwards, and the smallest ones downwards, by now a familiar phenomenon. Some authors aggregate stocks into portfolios, on the assumption that betas can be estimated more accurately for portfolios. Typically, portfolios are formed by ranking stocks on the basis of their betas estimated over a given period, then portfolio betas are estimated over a later period. This ensures that betas vary across portfolios, but prevents portfolio beta estimates from being biased. What this procedure actually does is shrink beta estimates together.

Since shrinkage is the general answer to such problems, why not apply the technique of Section 1.3? As it turns out, there is a direct correspondence between shrinking sample eigenvalues when T and N both go to infinity, and shrinking beta estimates (or sample means) when T is fixed and N goes to infinity. Thus, the asymptotic linear shrinkage developed in Sections 1.3.1-1.3.2 can be applied to betas too. However, linear shrinkage has no impact on t-statistics of regression slopes: it only changes the intercept. In other words, if the bias of betas is nearly linear, then there is little reason to fix it. For this reason,

we do not elaborate on this point here, and work with unadjusted betas. This more naive approach is less arbitrary than forming portfolios, and — if anything — makes it harder to find a significant relationship between returns and betas.

Previous OLS regressions of returns on betas found a positive slope, but with some serious limitations. First, it is not always statistically significant. Second, Tinic and West (1984) show that the return-beta relationship weakens substantially if the month of January is excluded from the period. Also, Lakonishok and Shapiro (1986) find that it disappears if size is included in the regression. Finally, Fama and French (1992) report that it flattens out over the period 1963-1992.

Using the same database as these authors, we reproduce their OLS results in Table E.3. The t-statistic for significance of the slope of returns on betas is 1.03 over the full period 1936-1992. It goes down to -0.33 if January is excluded, to 0.17 if size is included, and to 0.60 over 1963-1992. Actual results may differ somewhat from previously published ones, but the conclusions are identical.

Now, we change only one step: instead of using the structured estimator $\bar{\Sigma} = \hat{m}I$ for portfolio selection, we use a shrinkage estimator. This corresponds to upgrading from OLS to GLS. In Table E.4, we report the results obtained with the shrinkage estimator corresponding to the single index model, since this is the best-known structured estimator among the ones in Appendix B. The t-statistic for significance of the slope of returns on betas is now 1.91 over the full period 1936-1992. It is statistically significant at the 5% level against a one-sided alternative. It only goes down to 1.62 if January is excluded, to 1.44 if size is included, and to 1.16 over 1963-1992. This relationship is much more robust than under OLS.¹⁴

The change comes from two sources: standard deviations go down, because GLS is more efficient than OLS, and slope estimates go up. Kandel and Stambaugh (1994) explain theoretically why this should be anticipated. They show that OLS slope estimates can be more sensitive than GLS to misspecification of the market proxy used to estimate betas.

¹⁴The interested reader can find results for other asymptotic shrinkage estimators in Table E.5. All slope estimates are positive. The results that we choose to comment are neither the weakest nor the strongest, and are close to another structured estimator's results. We believe that they are the most credible.

In conclusion, the first predictive GLS cross-sectional regression of stock returns on betas, conducted thanks to the asymptotic linear shrinkage estimator of the covariance matrix developed in Section 1.3, finds a more significant and robust positive relationship between returns and betas than similar OLS regressions do. The relationship is not as strong as theory suggests, but this is hardly surprising given the error of beta estimates. Predictive GLS regressions support the existence of a positive linear relationship between returns and betas.

1.5 Conclusion

Directions for future research include using the spectral theory of large-dimensional random matrices to test for the number of factors in the APT; translating asymptotic shrinkage techniques to beta estimation; searching for the best frequency at which to sample stock returns for covariance matrix estimation; accounting for some type of Autoregressive Conditional Heteroskedasticity (ARCH) effects; bringing improved covariance matrix estimators to other areas of empirical stock market finance such as event studies.

In this chapter, we demonstrate the importance of a seldom-used framework for covariance matrix estimation: letting the number of variables and the number of observations go to infinity together. This framework is particularly well-suited for stock returns data, because the number of stocks traded in the stock market is at least of the same order of magnitude as the number of time periods. The covariance matrix of stock returns is important because it is a necessary input into portfolio selection, a central method in stock market finance.

We show that, in this framework, the sample covariance matrix is not well-behaved, especially through its eigenvalues. This work has potential implications for tests of the number of factors in the APT based on sample covariance matrix eigenvalues. We also show that it is easy to improve over the sample covariance matrix by shrinking its eigenvalues together in an asymptotically optimal way. In particular, this yields the first rotation-invariant estimator of the inverse of the covariance matrix to retain some theoretical motivation when variables outnumber observations. Generalizations provide attractive asymptotic extensions

to familiar finite sample Bayesian and decision theory methods.

Monte-Carlo simulations reveal that peak asymptotic performance is attained as soon as the number of observations and the number of variables become of order 20 to 40. The asymptotic shrinkage estimator has lower mean squared error than the sample covariance matrix in all situations simulated. It compares favorably overall in terms of mean squared error with existing finite sample estimators. The asymptotic shrinkage estimator has the potential to replace the sample covariance matrix as an all-purpose estimator.

More importantly for Finance, this asymptotic shrinkage technology helps portfolio selection minimize variance, as tests on historical data show. It opens new investment opportunities: equity portfolios with lower risk than was previously possible. These opportunities, however, are only slightly more attractive than existing ones because lower risk is penalized by lower return. In a related investigation of the risk-return tradeoff, the improved covariance matrix estimator is used to perform the first predictive GLS cross-sectional regression of returns on betas. This test concludes that the positive relationship between returns and betas predicted by the CAPM is statistically significant and robust, in stark contrast with tests based on less efficient OLS regressions.

Chapter 2

The δ -Arbitrage Pricing Theory

This chapter starts from three observations about the Arbitrage Pricing Theory (APT). First, the APT's assumption of no limiting arbitrage can be strengthened to characterize how closely beta pricing is approximated in a finite economy. Second, the stock market's risk structure does not show a clear-cut frontier that would validate the theoretical distinction between factors and residuals. Third, accounting for estimation error of factor risk premia, the optimal number of factors is determined by a trade-off between accuracy and parsimony in the beta pricing equation.

Based on these observations, I build a flexible and realistic model of the trade-off between risk and return in the stock market that capitalizes on previous research and opens interesting avenues for empirical work. The key assumption is to rule out the existence of δ -arbitrage opportunities, defined as portfolios whose Sharpe measure exceeds the predetermined level δ .

2.1 Introduction

The Arbitrage Pricing Theory (APT) yields an approximate beta pricing equation. However, its traditional form does not give researchers any way to check how closely the approximation holds in their economy. It is because the definition of limiting arbitrage (Ross, 1976; Huberman, 1982) is strong enough to obtain that the mean squared error on the beta pricing equation vanishes, but not strong enough to obtain the *rate* at which it vanishes. I propose

modifying the APT by ruling out δ -arbitrage opportunities, defined as portfolios whose Sharpe measure exceeds a certain level δ chosen by the researcher.¹ Realistic choices for δ are discussed. Ruling out δ -arbitrage is stronger than the APT's traditional assumption, but generally weaker than the Capital Asset Pricing Model's (CAPM). I give a simple illustrative model where ruling out δ -arbitrage follows from more primitive economic arguments. Stylized facts about the behavior of investors generally support this assumption.

To emphasize that the value of δ must be specified upfront, I call the strengthened model: δ -Arbitrage Pricing Theory. The δ -APT gives researchers a way to check how closely the approximate beta pricing equation holds in their economy. Unlike the APT, it still works fine if this error is not negligible, which may very well be the case in practice. This model admits as limiting cases: the APT with strict factor structure, the APT with approximate factor structure, the exact APT with noiseless residuals, and the CAPM. Some critics have claimed that the APT lacks economic content, but the δ -APT cannot be so criticized.

In the δ -APT, as opposed to the APT, it is not Nature but each researcher who chooses the factor structure. This is a desirable feature since researchers do disagree on the number and on the identity of "the" factors. Ignoring estimation error, it is always beneficial to project exogenous factors onto the space of portfolio returns. Optimal factors are returns on portfolios whose weights are eigenvectors corresponding to the top eigenvalues of the covariance matrix of stock returns.

I also examine the impact of estimation error. There is reason to believe that standard statistical techniques cannot distinguish factors from residuals in a predictive sense. In particular, it seems that the maximum residual eigenvalue is not as small as usually thought. Finally, accounting for estimation error in factor risk premia, the optimal number of factors is determined by a trade-off between accuracy and parsimony in the beta pricing equation.

In summary, this chapter proposes a more flexible and more realistic model of the familiar trade-off between risk and return. It opens new and promising directions for future empirical research.

Section 2.2 reviews the Arbitrage Pricing Theory. Section 2.3 defines δ -arbitrage.

¹The Sharpe measure of a portfolio is the ratio of the expectation to the standard deviation of its return in excess of the riskfree rate.

Section 2.4 presents the δ -Arbitrage Pricing Theory. Section 2.5 discusses the choice of factors. Section 2.6 evaluates the impact of estimation error. Section 2.7 concludes.

2.2 Arbitrage Pricing Theory

The APT is reviewed and a fundamental limitation is identified.

2.2.1 Review of the APT

Following Shanken (1985), I consider that equilibrium derivations of the APT are a breed apart. They are outside the scope of this chapter.

Asset returns are generated by the factor model:

$$\tilde{r}_i = \beta_{i1}\tilde{f}_1 + \dots + \beta_{iK}\tilde{f}_K + \tilde{e}_i, \quad (2.1)$$

where \tilde{r}_i is the return on the i^{th} asset ($i = 1, \dots, N$) in excess of the riskfree rate, $\tilde{f}_1, \dots, \tilde{f}_K$ are factors, $\beta_{i1}, \dots, \beta_{iK}$ are factor loadings (also called betas), and \tilde{e}_i is the residual. By definition, residuals are uncorrelated with factors.

Let $\bar{\lambda}$ denote the largest eigenvalue of the covariance matrix of the residuals $(\tilde{e}_i)_{i=1, \dots, N}$. Intuitively, if $\bar{\lambda}$ is not large, then residuals do not explain much of the risk: factors explain a lot of it.

Let $\bar{\delta}$ denote the largest Sharpe measure in the economy. The Sharpe measure of a portfolio is defined as the ratio of the expectation to the standard deviation of its return in excess of the riskfree rate. Intuitively, if $\bar{\delta}$ is not large, then there is a tight relationship between risk and return.

In the APT, the beta pricing equation:

$$E[\tilde{r}_i] \approx \beta_{i1}\tau_1 + \dots + \beta_{iK}\tau_K \quad (2.2)$$

holds approximately for an appropriate choice of factor risk premia τ_1, \dots, τ_K . It means that betas drive most of the expected return in the economy. The mean squared approximation

error in Equation (2.2) is:

$$\varepsilon^2 = \min_{\tau_1, \dots, \tau_K} \frac{1}{N} \sum_{i=1}^N (E[\tilde{r}_i] - \beta_{i1}\tau_1 - \dots - \beta_{iK}\tau_K)^2. \quad (2.3)$$

If ε can be neglected, then beta pricing is a reliable approximation.

The APT is the following result.²

Theorem 11 (APT) *If the number of assets N goes to infinity while the maximum residual eigenvalue $\bar{\lambda}$ and the maximum Sharpe measure $\bar{\delta}$ remain bounded, then the mean squared approximation error ε^2 of the beta pricing equation (2.2) vanishes.*

This succinct review masks some of the subtler points involved in allowing the number of assets to go to infinity. See Ross (1976) or Chamberlain and Rothschild (1983) (hereafter CR) for more detailed presentations of the APT.

2.2.2 Limitation of the APT

The APT is an asymptotic theory. Its implications are for N going to infinity: it says nothing about a finite economy. In practice, only a finite number of stocks are traded, so some critics wonder whether the APT is at all useful.

Superficially, one can write off their concern as an elementary misunderstanding of asymptotics: if there are enough stocks in the economy for the APT's beta pricing equation to be “approximately” true, then the APT is useful. At a deeper level, however, these critics may just be asking for evidence that the APT is “approximately” true.

In other words, for the APT to have practical relevance, it cannot just state that the mean squared approximation error ε^2 on the beta pricing equation (2.2) converges to zero. The APT must also say *at what rate* ε converges to zero, so that critics can check how close the approximation is in reality.

Let us re-examine CR's proof of convergence to determine the convergence rate. The crucial equation in the proof is:

$$\varepsilon^2 \leq \frac{\bar{\lambda} \bar{\delta}^2}{N}. \quad (2.4)$$

²Proofs are in the Appendix.

Therefore critics of the APT should plug in reasonable values for $\bar{\lambda}$, $\bar{\delta}$ and N , then check whether these values imply a tight enough error bound.

Getting the number of assets N is obviously no problem. The maximum residual eigenvalue $\bar{\lambda}$ can be directly estimated from the data, in principle. There are some difficulties: the location of the residual space is not clear (see Section 2.5), and estimation error may be severe (see Section 2.6). But at least it is theoretically feasible.

Now we are facing the main problem: the APT says nothing about the maximum Sharpe measure $\bar{\delta}$ (except that it is finite). Since the APT does not specify $\bar{\delta}$, it prevents us from checking how closely the beta pricing equation (2.2) holds in our economy.

An asymptotic approximation without its convergence rate is useless in a finite economy.

This is a severe, fundamental limitation of the APT. It was already noted by Shanken (1982). Thankfully, the way to get around it has been known intuitively since the origin of the APT. I formalize it below.

2.3 δ -Arbitrage

Alternative definitions of arbitrage are reviewed and a new one is proposed. Its economic justification is discussed qualitatively and illustrated by a quantitative model. Some empirical evidence is presented and values are recommended for the key parameter.

2.3.1 Definitions of Arbitrage

Ross's original intention was to demonstrate that the trade-off between risk and return did not require the CAPM's strong equilibrium assumption. More generally, it is characteristic of Finance (as opposed to Economics) to prove results under an assumption weaker than equilibrium. This assumption is usually that no arbitrage opportunities exist.

An arbitrage opportunity can be defined as a riskfree investment that earns more than the riskfree rate. Its return in excess of the riskfree rate has strictly positive expectation but zero standard deviation. In the mean-variance world characteristic of stock market studies, an arbitrage opportunity is an investment with infinite Sharpe measure.

Ruling out arbitrage opportunities is a very weak restriction on agent behavior, hence

a very credible one. However, it is so weak that it has no implication for the economy in Section 2.2.1.³ Therefore it has to be strengthened. The APT rules out *limiting* arbitrage opportunities.

A limiting arbitrage opportunity is a sequence of portfolios whose Sharpe measures go to infinity. Ruling out limiting arbitrage is equivalent to assuming that the maximum Sharpe measure $\bar{\delta}$ is finite. Exactly how finite is not said. This is the root of the problem outlined above.

The solution was already indicated by Ross (1976): in a back-of-the-envelope calculation, he assumes that the maximum squared Sharpe measure $\bar{\delta}^2$ is less than twice the squared Sharpe measure of a market index. It is a small step from saying that $\bar{\delta}$ is finite but unspecified to saying that $\bar{\delta}$ is finite *and specifying it*. Yet this step generates a new notion of arbitrage and fixes the fundamental problem of the APT.

More generally, any empiricist using the approximate beta pricing equation makes the leap of faith that it is accurate enough for practical purposes. Implicit here is the assumption that the maximum Sharpe measure in the economy is not large enough to undermine the accuracy of beta pricing. Therefore the model formally developed below is nothing more than the one that proponents of the APT had in the back of their minds all along.

2.3.2 δ -Arbitrage

I give a definition of arbitrage that differs from the APT's limiting arbitrage.

Definition 3 A δ -arbitrage opportunity is a portfolio whose Sharpe measure is strictly above δ .

The value of δ is *specified a priori* by the economist. Different economists may legitimately disagree on the value of δ . I will talk about δ -arbitrage in general terms, so that everyone has a chance to replace δ by the value that they deem appropriate. I will discuss what values I deem appropriate in Section 2.3.6.

The APT rules out the existence of limiting arbitrage opportunities. I need a stronger assumption.

³Except if the covariance matrix of stock returns is singular, a case treated by Ross (1978) but not realistic.

Assumption 4 *There are no δ -arbitrage opportunities.*

An asset pricing theory based on Assumption 4 can be seen as a whole family of models indexed by δ .

If $\delta_1 < \delta_2$ then ruling out δ_1 -arbitrage implies ruling out δ_2 -arbitrage. The strength of Assumption 4 decreases in δ . Therefore the model indexed by δ_2 is embedded into the model indexed by δ_1 . The strength of the conclusions of the models decreases in δ .

In the limit as δ goes to infinity, Assumption 4 converges to the assumption of no limiting arbitrage. Therefore the union of all the models in the family is nothing else than the APT itself (if the appropriate auxiliary assumptions are added).

In the opposite direction, δ cannot be set below the Sharpe measure of any particular portfolio in the economy. For example, δ cannot be set below the Sharpe measure δ_M of the market portfolio. In the limit if $\delta = \delta_M$, then Assumption 4 states that the market portfolio is mean-variance efficient. Therefore the intersection of all the models in the family is nothing else than the CAPM itself (again if the appropriate auxiliary assumptions are added).

In conclusion, ruling out δ -arbitrage is a flexible way to bridge the gap between two important economic assumptions: no-arbitrage and equilibrium. The corresponding asset pricing theory bridges the gap between the APT and the CAPM. This formalism lets economists modulate at will the strength of their restriction on investor behavior, and reach a conclusion of accordingly variable strength. The next section examines in more detail the economic justification of Assumption 4.

2.3.3 Economic Justification

The Sharpe measure is an economically justified criterion for the attractiveness of a portfolio if agents only care about risk (negatively) and return (positively), and if riskfree borrowing and lending are available at the same rate. In such a world, mean-variance efficiency is attained by investing in the safe asset and in the portfolio of risky assets with the highest Sharpe measure. As a consequence, every agent holds a scalar multiple of the same portfolio of risky assets: the one with maximum Sharpe measure. By aggregation, the

market portfolio must have maximum Sharpe measure, as the well-known proof of the CAPM goes.

The disturbing point is that, in reality, all agents do not hold the same portfolio. Surely then they must care about (personal) objectives other than risk and return. In this case, there is little reason to believe that the market portfolio has maximum Sharpe measure. The portfolio with maximum Sharpe measure strictly dominates the market portfolio (in combination with the riskless asset), yet some agents choose not to hold it for personal reasons.

How high can the maximum Sharpe measure be before the corresponding portfolio becomes overpoweringly attractive? Arbitrarily high? Certainly not! There has to be a limit on the intensity of the personal objectives that distract agents away from maximizing their Sharpe measure. If there is, it implies a limit on the maximum Sharpe measure that can prevail in equilibrium. If an investment opportunity's Sharpe measure is too attractive, distractions will not be sufficient to keep agents from arbitraging it away. The less agents care about other objectives, the lower the maximum Sharpe measure. If agents do not care at all about non-mean-variance objectives, then the maximum Sharpe measure is as low as the market portfolio's (CAPM).

Therefore, qualitatively, δ -arbitrage is ruled out if there is an upper limit on the intensity of personal objectives that distract agents away from looking at Sharpe measures. It is certainly not the only or even the best way to transform a limit of this kind into a quantitative restriction. Particular specifications of personal objectives may lead to different restrictions. But one important advantage of Assumption 4 is that it makes minimal assumptions on the nature of these objectives, about which so little is known.

In summary, ruling out δ -arbitrage is a general restriction with powerful implications that characterizes the intensity, rather than the nature, of the reasons why different investors hold different portfolios.

2.3.4 Non-Marketable Assets

The above discussion was purely qualitative, but it can also be pursued quantitatively. Consider a specific example of what might distract agents away from their Sharpe measures. Roll (1977) argues that the existence of assets outside the stock market makes it difficult to conduct asset pricing in the stock market in isolation. The risk of non-marketable assets such as human capital may influence the pricing of marketable assets such as stocks.

In the following example, a bound on the riskiness of non-marketable assets implies a bound on the maximum Sharpe measure in the stock market.

Theorem 12 (Equilibrium with Non-Marketable Assets) *Assume that agents have constant absolute risk aversion and that future asset values are normally distributed. Let σ_M^2 (respectively σ_{NM}^2) denote the variance of the future dollar value of all marketable (resp. non-marketable) assets in the economy. If the correlation between marketable and non-marketable assets is non-negative then, in equilibrium, the maximum Sharpe measure in the market $\bar{\delta}$ verifies:*

$$\bar{\delta}^2 \leq \delta_M^2 \left(1 + \frac{\sigma_{NM}^2}{\sigma_M^2} \right), \quad (2.5)$$

where δ_M is the Sharpe measure of the market portfolio.

The notion of equilibrium here differs from the CAPM's because agents pursue objectives besides mean-variance efficiency of marketable assets.

The assumption that the aggregate value of human capital has non-negative correlation with the aggregate value of companies in the stock market is at least plausible. As an illustration, if the variance of the future value of human capital is no greater than the variance of the future value of the market portfolio, then the maximum squared Sharpe measure $\bar{\delta}^2$ cannot exceed the bound $\delta^2 = 2\delta_M^2$.

2.3.5 Empirical Evidence

Peter Lynch is reputed for his outstanding tenure as manager of the Fidelity Magellan mutual fund. From 1977 to 1990, his fund had about one and a half times the Sharpe measure of a market index. In our model, we expect investors to flock massively to such

investment opportunities, thereby disrupting prices. Did they flock to Fidelity Magellan? Yes: funds under management grew from \$20 million in 1977 to \$13 billion in 1990! Did they disrupt prices? It is not as easy to assess it directly, if only because Peter Lynch followed a dynamic strategy with frequent rebalancing. However, the sheer size of his fund severely limited the fraction invested in any given stock. It became more time-consuming for him over the years to amass the private information necessary to pick enough winners to sustain his performance. For compensation, he charged steep commission fees that deflated his shareholders' Sharpe measure. Eventually, his commission did not compensate him enough for the time spent at work and he resigned, taking his talent away from the fund. This story is anecdotal evidence that Assumption 4 is a sensible description of the world.

MacKinlay (1993) and Daniel and Titman (1995) argue that non-risk characteristics of stocks are priced. Size and book-to-market are often cited as examples. This is interesting evidence that *appears* to violate Assumption 4. As these authors point out, if the phenomenon is real and persists, then approximate arbitrage opportunities will exist. But their results do not, alone, constitute evidence contrary to Assumption 4. Contrary evidence would be if mutual funds exploiting this approximate arbitrage were set up, convinced the public that they have high Sharpe measures and, in spite of that, remained terribly unpopular! Assumption 4 restricts investor behavior, so investor behavior alone can disprove it. I doubt that high Sharpe measure mutual funds loading on priced non-risk characteristics can remain unpopular for long, if they perform that well. And if they do become popular, they will find it harder to maintain their performance. For example, the size effect seems to have disappeared a few years after its discovery. I am confident that any apparent violation of Assumption 4, as soon as it is established beyond a doubt, must quickly go away.

2.3.6 Choice of δ

Ross's (1976, p. 354) back-of-the-envelope calculation, Section 2.3.4, the example of Fidelity Magellan and MacKinlay's (1993, Section 5.2) risk-based alternative specification all suggest that a reasonable choice for δ^2 might be in the neighborhood of twice the squared Sharpe measure δ_M^2 of the market portfolio. Since the spirit of the δ -APT requires that I

recommend a specific value for δ , at this stage I choose $\delta = \sqrt{2}\delta_M$.

Of course, disagreement is legitimate. CAPM theorists insist that δ must be equal to δ_M . Fama and French (1992) believes that only a Sharpe measure of the order of 150 δ_M comes close to the intuitive notion of arbitrage. Both positions are somewhat extreme. The question is: as an investor, would you go out of your way to make an investment earning the same expected return as the market with half the variance? I believe that most investors would answer positively.

2.4 δ -Arbitrage Pricing Theory

The assumption of no δ -arbitrage generates a modified version of the APT. Its relation to existing asset pricing theories, economic contents and testability are examined in this section.

2.4.1 Formulation

Recall the universe described in Section 2.2.1 before Theorem 11. Its key elements are the number of assets N , the maximum residual eigenvalue $\bar{\lambda}$, and the maximum Sharpe measure $\bar{\delta}$.

Factors can be projected onto asset returns:

$$\tilde{f}_k = m_{k1}\tilde{r}_1 + \dots + m_{kN}\tilde{r}_N + \tilde{\eta}_k, \quad (2.6)$$

where the projection residual $\tilde{\eta}_k$ is uncorrelated with asset returns. The coefficients m_{ki} are weights of factor-mimicking portfolios. The projection residual $\tilde{\eta}_k$ has zero variance if and only if the k^{th} factor is spanned by asset returns. Generally, this is not the case. Let δ_F denote the maximum Sharpe measure among the portfolios that are spanned by the K factor-mimicking portfolios.

This universe is entirely non-restrictive. In particular, the number of assets N can be anything and does *not* go to infinity. No economic assumption has been made (yet). In this finite economy, the δ -Arbitrage Pricing Theory is the following result.

Theorem 13 (δ -APT) *With the above notation, if Assumption 4 holds, then the mean squared error ε^2 on the approximate beta pricing equation (2.2) verifies:*

$$\varepsilon^2 \leq \frac{\bar{\lambda} (\delta^2 - \delta_F^2)}{N}. \quad (2.7)$$

Intuitively: If factors explain most of the risk in the economy ($\bar{\lambda}$ small), if approximate arbitrage opportunities are ruled out (δ not large), and if many assets trade (N large), then betas drive most of the expected return in the economy (ε small).

Theorem 13 captures the intuition of Ross's (1976) and CR's Arbitrage Pricing Theory. In addition, thanks to the strength of Assumption 4, Equation (2.7) shows how closely beta pricing holds in any given finite economy.

The presence of δ_F makes the bound sharper than the one in Equation (2.4) copied from CR. Shanken (1982) proves a similar mathematical result, but does not place it within a formal economic theory.

The intuition behind the δ -APT is not new, since Ross (1976) already plugged a specific value into δ . The mathematics are not new either, since they are a minor rewriting of CR's result. Shanken (1992) even goes so far as to bring the intuition and the maths together to call for empirical work on the relationship between deviations from beta pricing and approximate arbitrage opportunities.

But this is the first time that the definitions of δ -arbitrage and δ -Arbitrage Pricing Theory are given. My contribution is to show that setting δ to a specific value is more than an eloquent illustration or an empirical convenience: it is an integral part of the theory that guarantees its relevance and was previously missing.

It could be said that the intuition for my work is not novel because it was implicitly assumed in every discussion of the practical relevance of the APT. I gladly admit that my only contribution is to make explicit an idea that was previously implicit. But this contribution changes the focus, the interpretation and the span of the theory. The major advantage is that important empirical issues that were previously outside the APT can now be treated inside the δ -APT.

2.4.2 Relation to Existing Asset Pricing Theories

Equation (2.7) can be taken to the limit in various ways.

The pricing error ε vanishes as N goes to infinity. This is the traditional APT result.⁴ However Assumption 4 buys us more than that: we also know at what *rate* the error vanishes. If different people disagree on δ , then they disagree on the rate, even though they still agree that ε vanishes. The only thing that is implied by every choice of δ is that ε vanishes, which is the traditional APT result. In practice, it seems likely that the research community will at least reach a consensus on a (finite) range for δ , which would avoid going back to the traditional APT.

The original objective of δ -arbitrage was to allow asymptotic analysis of the rate of convergence to beta pricing. This objective has been fulfilled. However, since the bound (2.7) holds not just asymptotically but also in any finite economy, there is no need to invoke the original asymptotic formalism anymore. Again, the δ -APT is *not* an asymptotic theory.

Notice that the pricing error ε vanishes as the maximum residual eigenvalue $\bar{\lambda}$ goes to zero. This corresponds to the exact pricing result in an economy with noiseless residuals proven by Ross (1978).

Also, setting the number of factors to one and taking the single factor as the return on the market portfolio ensures that δ_F is equal to the Sharpe measure δ_M of the market portfolio. Then, if economists stop allowing investors to pursue other objectives than the Sharpe measure, the δ from Assumption 4 goes to δ_M , therefore $\delta - \delta_F$ vanishes, and so does the error in the beta pricing equation. This is the CAPM result.

Previous versions of the APT and the CAPM were designed to give conditions under which the error ε is negligible. By contrast, the δ -APT works with any value of ε , negligible or not. This could prove a decisive advantage if empirical results show that ε is *not* negligible, which they may very well do. In this case, it would be preferable to take ε into account explicitly and know its order of magnitude, as in the δ -APT, rather than neglect it because some other version of the theory cannot accommodate it.

⁴Shanken (1992, Section II) reports this connection.

In conclusion, the δ -Arbitrage Pricing Theory is flexible enough to admit existing versions of the APT and the CAPM as limit cases. But the drawback of these limit cases is that they are, strictly speaking, unrealistic. A more realistic approach is to evaluate whether, together, the large number of assets in the economy, the low risk of residuals, and the limit on acceptable Sharpe measures *combine* to make a beta pricing equation accurate. One of the major contributions of this chapter is to develop a theoretical framework where this approach is permitted.

2.4.3 Economic Contents

One of Roll's (1977) contributions is to dissect the CAPM into:

1. The mathematical equivalence between the mean-variance efficiency of the market portfolio and the equation linking expected returns to betas;
2. The economic assumption that the market portfolio is mean-variance efficient.

I claim that the δ -APT has the same structure. It can be decomposed into:

1. The mathematical equivalence in Theorem 13 between bounding the maximum Sharpe measure and bounding the error of the beta pricing equation (2.2);
2. The economic assumption that places a specific upper bound on the maximum Sharpe measure.

As Shanken (1992) points out, Theorem 13 is a mathematical tautology. Shanken (1982) concludes from a repackaging argument that the APT is an empty theory with zero economic content. The above interpretation makes it clear that these are not valid criticisms of the δ -APT. The δ -APT asks economists to take a stand on what Sharpe measures are reasonable. This in turn translates into a beta pricing equation, which is more or less accurate, depending on how assertive economists are about the maximum Sharpe measure. In the limit, by setting the maximum Sharpe measure equal to the Sharpe measure of the market portfolio, the δ -APT can be turned into the CAPM, a theory of (perhaps excessive but) surely undisputed economic content. To sum up, the economic content in the δ -APT is substantial and it is... δ !

2.4.4 Testability

To push this reasoning one step further, can δ be mechanically estimated, bypassing the need for economic judgment? Strictly speaking, only $\bar{\delta}$ could be estimated, and the economist (not the statistician) would bear the ultimate responsibility for plugging the estimate into δ . But the question is still interesting.

Such estimation would require knowledge of expected stock returns. If expected returns are known, then the beta pricing equation is not useful. Normally, expected stock returns are not known, and Equation (2.2) serves to compute them from betas and factor risk premia, which are known (or at least estimated with less error). If everything is known, then verifying that the mathematical tautology (2.7) holds is a waste of time.

Conceivably, $\bar{\delta}$ could be estimated over the past, the estimate plugged into δ , and Equation (2.2) used to make future investments. Then one could hesitate between using the historical estimate of $\bar{\delta}$ and using the value of δ suggested by economic judgment (see Section 2.3.6). To some extent, economic judgment relies on accumulated knowledge about agent behavior and, as such, is just another historical estimate! Conversely, the belief that investor behavior will not change from the estimation period to the investment period (even if δ -arbitrage opportunities are publicly reported in the meantime) is itself an economic judgment...

In this particular case, it would come as no surprise if the *ex post* knowledge that, say, size and book-to-market commanded non-risk based premia over 1963-1990, generated an estimate of $\bar{\delta}$ higher than the δ that I recommend in Section 2.3.6 (see MacKinlay, 1993). Which one to believe: the high estimate or the low recommendation? It all hinges on whether non-risk based premia can persist. Over long horizons, I would put more faith in the efficacy of δ -arbitrageurs than in the perpetual reappearance of free lunches.

This discussion shows that economic judgment is essential in the choice of δ and cannot be totally replaced by an automatic estimation procedure. After all, it is the *ex ante* δ that we are concerned with, and economic thinking sheds more light on *ex ante* parameters than *ex post* data sometimes do.

To sum up, Assumption 4 is testable: for example, MacKinlay (1993) tests it. However,

the conclusions of such tests has to be carefully moderated by economic judgment.

Should convincing empirical evidence accumulate against the version of Assumption 4 recommended in Section 2.3.6, the whole theory of δ -APT does not fall apart. It would suffice to raise the value of δ until it conforms with this evidence, at the cost of weakening the conclusions of the model.

2.5 Choice of Factors

The δ -APT provides a convenient framework to study one of the most controversial problems about empirical applications of the APT: the choice of factors. Contrary to the APT, the δ -APT does not assume that the factors are uniquely determined by the risk structure of the stock market. Factors could be anything. Of course, not all sets of factors are created equal. The number and the identity of the factors matter through the maximum residual eigenvalue $\bar{\lambda}$. The rule is: the lower the $\bar{\lambda}$, the more accurate beta pricing is.

2.5.1 Factors vs. Residuals

A distinction central to the APT separates factors from residuals. In Ross's (1976) APT with strict factor structure, factors are pervasive while residuals are idiosyncratic. In Ross's (1978) noiseless APT, factors are risky while residuals are riskless. In Chamberlain and Rothschild's (1983) APT with approximate factor structure, factors correspond to eigenvalues going to infinity and residuals to bounded eigenvalues.

I argue that data do not support the theoretical distinction between factors and residuals. True, market risk is an important source of risk in the stock market and, as such, fits well the theoretical definition of a factor. Also true, some sources of risk have such minute influence on stock returns that they can be called residuals. But there is almost a continuum of sources of risk filling the gap between these two extremes. Observed risk structure does not show a clean frontier between two groups that would embody the theoretical distinction between factors and residuals.

Figure F-10 plots the sorted eigenvalues of the covariance matrix of returns on stocks trading on New York Stock Exchange (NYSE) from 1988 to 1993. The first eigenvalue is

much larger than the last ones, but no evident gap jumps to the eye.

Notice that the difference between consecutive eigenvalues fades away as eigenvalues become smaller. If we simply want the largest possible gap, we must call the first eigenvalue a factor and all the other ones residuals; but then residuals would explain more than 90% of the total risk in the stock market, hardly an intuitive feature! Now, if we want to be able to neglect the variance explained by residuals, we must have a very large number of factors; but then the difference between the smallest factor and the largest residual becomes uncomfortably small: did we get exactly the right number of factors, or one too few? or one too many? can we tell?

This does not sound like a satisfactory way to determine the number of factors in the APT, yet similar reasonings have been made in countless empirical studies of the question. We should accept that Nature does not separate factors from residuals: researchers have to do it themselves. Put another way, different researchers can legitimately disagree on what “the” factors are. This may explain why researchers do in fact disagree.

The fact that researchers are free to choose the factors they want does not preclude some choices from being better than others. This matter is discussed below and again in Section 2.6.3.

2.5.2 Exogenous Factors

What can determine the choice of factors? On the one hand, the more risk factors explain, the lower $\bar{\lambda}$ is, and the more binding Equation (2.7) is. On the other hand, the higher the number of factors, the more trivial pricing equation (2.2) becomes: explaining 60 assets with 50 factors would be plain silly. Therefore we should choose factors to reduce $\bar{\lambda}$ without increasing K .

Let us start from an exogenously specified set of factors $\tilde{f}_1, \dots, \tilde{f}_K$. In general, they are not spanned by asset returns. Remember that they can be projected onto the space of portfolio returns, as in Equation (2.6). The projection of the k^{th} factor is the return on the k^{th} mimicking portfolio. The following result shows that it is advantageous to replace factors by mimicking portfolio returns.

Theorem 14 *Replacing factors by their projections onto asset returns reduces the maximum eigenvalue of the covariance matrix of residuals.*

This is an easy way to give more bite to the δ -APT. In practice, this projection requires the inverse of the covariance matrix of stock returns. If the number of stocks in the cross-section exceeds the number of time-series observations, which is often the case, then the usual estimator of the covariance matrix is not invertible. The asymptotic shrinkage technique devised in the first chapter of this thesis gets around this problem and obtains a covariance matrix estimator that is guaranteed to be always invertible. This estimator can be used to project factors onto stock returns as in Theorem 14.

2.5.3 Covariance Matrix Eigenvectors

Holding the number of factors K fixed, how do we choose K factors to minimize $\bar{\lambda}$? By Theorem 14, optimal factors are returns on carefully chosen portfolios. But what portfolios? The following theorem answers.

Theorem 15 *Let the k^{th} factor be the return on the portfolio whose weight vector is the eigenvector of the covariance matrix of asset returns corresponding to the k^{th} largest eigenvalue. This choice of factors minimizes the largest residual eigenvalue $\bar{\lambda}$.*

The way to give the most bite to the δ -APT is to choose the factors associated with the first K eigenvectors of the covariance matrix of asset returns. Researchers who rely on exogenous factors (see e.g. Chen, Roll and Ross, 1986) should be aware that their residuals might be very risky, which would in turn compromise the accuracy of their beta pricing equation.

2.6 Estimation Error

Much of the work on the APT assumes that the mean vector and covariance matrix of stock returns are known exactly. In reality, they are estimated with error. How does this affect beta pricing? The δ -APT gives a convenient framework to start exploring this important question.

2.6.1 Forecasting Residual Space

First, the effect on δ -arbitrageurs. These are the investors who locate the factor space and the residual space, determine whether some residuals command high returns and, if so, invest in them.⁵ It is the action of δ -arbitrageurs that guarantees that residuals cannot consistently earn high returns. These agents focus on residuals rather than factors because it would be unlikely for any portfolio to earn returns so high that they exceed fair compensation for factor risk.

How does estimation error prevent δ -arbitrageurs from enforcing the δ -APT? They have to estimate the factors and the residual space in one period, then possibly invest in the most attractive portfolio in the residual space over the next period. At a minimum, they have to forecast where the factor space and the residual space will lie. If their forecast is always wrong because of estimation error, they will soon take it into account and possibly reduce their involvement in δ -arbitrage.

For example, residuals should be uncorrelated with factors. But some portfolios in the estimated residual space may turn out to have high correlations with factors over the investment period. If δ -arbitrageurs cannot weed out factor risk from their residuals, they may stop trying to take advantage of residual returns for fear of unintentionally bearing factor risk. Evidence that standard statistical techniques cannot discriminate factors from residuals in a predictive way would cast serious doubts on the existence of δ -arbitrageurs, and on the δ -APT altogether. This is an empirical question that will be explored in future research.

2.6.2 Maximum Residual Eigenvalue

Ideally, all factor eigenvalues should exceed all residual eigenvalues (see Theorem 15). However, if the eigenstructure is estimated with error from historical data, this may not be the case over the investment period. I conjecture that standard statistical techniques may not be accurate enough to prevent the largest residual eigenvalue from exceeding the smallest

⁵In the spirit of this chapter, a portfolio earns “high” return if it earns much more than fair compensation for its risk, i.e. its Sharpe measure exceeds δ .

factor eigenvalue.

Empirical evidence supporting this conjecture would be the last nail in the coffin for the idea that Nature separates factor risk from residual risk. In my opinion, a more realistic position is that the risk structure of the stock market gives researchers the freedom to disagree about the number and the identity of the factors. The δ -APT, unlike the APT, perfectly accommodates this position.

Even though the δ -APT will suffer less than the APT if this conjecture is verified, it will still suffer somewhat. The reason is that the accuracy of the beta pricing equation depends on the maximum residual eigenvalue $\bar{\lambda}$ that will prevail *over the investment period*. I believe that this value of $\bar{\lambda}$ is rather high, possibly higher than the minimum factor eigenvalue, and certainly higher than the maximum residual eigenvalue that prevailed over the estimation period.

In summary, eigenstructure estimation error increases the error in the beta pricing equation. Evaluating the severity of this effect is an urgent direction for future empirical research.

2.6.3 Optimal Number of Factors

The δ -APT provides a convenient framework to determine the optimal number of factors. Obviously, the first objective in choosing factors is to minimize $\bar{\lambda}$, the largest residual eigenvalue. This can be accomplished by selecting more and more factors. In the limit, selecting as many factors as there are assets yields $\bar{\lambda} = 0$, therefore the beta pricing equation is perfectly accurate. However, this is hardly a reasonable choice.

There must be a way to penalize sets of factors that are not parsimonious. The basic idea of the APT is that knowledge of the expected return on a few combinations of assets (the factors) contains information about expected returns on all assets. The premise here is that we know the expected returns on a few factors but not on all stocks. A more modest statement is that estimating the mean returns on a few factors can be done more accurately than estimating the mean returns on all individual stocks.

There are several reasons why this might be the case. First, it is easier to estimate a

few parameters than many parameters. Second, factors are less risky than stocks because they are better diversified, so their mean return estimates have lower standard errors. Third, if factor risk premia are characteristic of the economy as a whole, they may change more slowly than individual stocks: companies are created and liquidated but the economy stays the same.

A realistic way to encourage parsimony in choosing the number of factors is to take into account explicitly estimation error of factor risk premia. Let us work in the universe of Sections 2.2.1 and 2.4.1. Assume that factors are chosen optimally as in Theorem 15. The k^{th} factor is the return on the portfolio whose weight vector is the eigenvector corresponding to the k^{th} largest eigenvalue of the covariance matrix of stock returns.

For convenience, I assume that the covariance structure is known exactly. This is rather optimistic in view of Sections 2.6.1 and 2.6.2, but second moments of stock returns are better known than first moments. One important reason is that, keeping the estimation period fixed, infinitely increasing the sampling frequency yields consistency of estimators of the second, but not the first, moments. For example, moving from the monthly to the daily database from the Center for Research in Security Prices (CRSP) over the same estimation period increases the precision of the sample covariance matrix but not of the sample mean vector of stock returns.⁶

Let T denote the number of observations available. I assume that they are independent and identically distributed (iid). The k^{th} factor risk premium τ_k is estimated by the sample mean $\hat{\tau}_k$ of the return on the k^{th} factor over the T observations. Note that the variance of the return on the k^{th} factor ($k = 1, \dots, K$) is the k^{th} largest eigenvalue λ_k of the covariance matrix of stock returns. Therefore the estimation error on the k^{th} factor risk premium is: $E[(\hat{\tau}_k - \tau_k)^2] = \lambda_k/T$. Also note that the largest residual eigenvalue $\bar{\lambda}$ is equal to λ_{K+1} . Theorem 13 can be modified to account for risk premium estimation.

Theorem 16 *With the above notation, if Assumption 4 holds, then the accuracy of the beta*

⁶To mitigate this, daily data need to be cleansed from spurious microstructure effects more carefully than monthly data do. But there are ways to do it

pricing equation with estimated factor risk premia is determined by:

$$E \left[\frac{1}{N} \sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \hat{\tau}_k \right)^2 \right] \leq \frac{\lambda_{K+1} (\delta^2 - \delta_F^2)}{N} + \frac{1}{NT} \sum_{k=1}^K \lambda_k. \quad (2.8)$$

The first term on the right hand side of Equation (2.8) is the familiar bound from Equation (2.7). The second term contains risk premium estimation error. As anticipated, the first term decreases in K and the second one increases in K .

Equation (2.8) shows how to determine the optimal number of factors in the δ -APT as a trade-off between accuracy ($\lambda_{K+1}(\delta^2 - \delta_F^2)/N$) and parsimony ($\sum_{k=1}^K \lambda_k/NT$) in beta pricing equation (2.2). Of course, this trade-off depends critically on the length of the estimation period T .

Theorem 16 is also relevant for the choice between exogenous factors and eigenvectors. With exogenous factors, residuals are more risky, but the estimation period can be longer: most scholars would agree that the assumption that the covariance structure of stock returns is stationary over twenty years is heroic, yet they would not hesitate to estimate the market risk premium over 1926-1993. It is an open empirical question whether this trade-off favors exogenous factors or eigenvectors.

2.7 Conclusion

The δ -Arbitrage Pricing Theory reconciles APT formalism with the interpretation of the APT prevailing throughout empirical work. It does so by strengthening the economic assumption from no limiting arbitrage to no δ -arbitrage.

Ruling out δ -arbitrage opportunities, defined as portfolios whose Sharpe measures exceed the level δ , follows from placing a limit on how intensely agents are distracted away from the pursuit of mean-variance efficiency. It seems that a consensus can be reached on a range of values appropriate for δ .

The δ -APT covers a family of models stretching from Ross's (1976) APT to the CAPM. It lets economists select the model that is compatible with what they know of investor behavior. The δ -APT compares favorably in realism with existing asset pricing theories.

Specifying the value of δ , which is the essential point of the theory, makes the δ -APT testable.

One of the most valuable aspects of the δ -APT is that it sheds new light on difficult implementation issues such as the number of factors, their identity, and estimation error. It forms a solid foundation on which to base future empirical work.

Chapter 3

Is Beta Pricing Accurate?

This is an empirical investigation of what the Arbitrage Pricing Theory (APT) can and cannot do. The APT states that betas with respect to the major factors of risk in the stock market explain expected returns on stocks, or else approximate arbitrage opportunities would exist. I find that the link between deviations from beta pricing and approximate arbitrage is weak in practice. Even if beta pricing makes an error of $\pm 3.5\%$ on every expected return (quoted on an annual basis), the maximum Sharpe ratio in the stock market need not be more than one and a half times the Sharpe ratio of a value-weighted index. Thus, large deviations from beta pricing are compatible with the absence of approximate arbitrage opportunities.

This result goes against the spirit of Ross's (1976) APT, yet it is obtained simply by applying the mathematics of the APT to historical stock returns data. Deviations from beta pricing come from riskiness of residuals and from estimation error on factor risk premia.

Section 3.1 describes the empirical strategy. Section 3.2 investigates the APT with a single exogenous factor proxying for the return on the market portfolio. Section 3.3 considers multiple endogenous factors corresponding to eigenvectors of the covariance matrix of stock returns. Section 3.4 concludes.

3.1 Empirical Strategy

This section reviews the theoretical foundations for beta pricing, outlines the overall objective of this study, enumerates the main issues, and describes the data.

3.1.1 Beta Pricing

The following assumptions are maintained throughout. There is one riskless bond and a finite number of risky stocks. Stock returns have finite second moments. No combination of stocks is riskless. Agents can buy or sell any amount of stocks and bond without frictions. In this world, the mean-variance efficiency of a portfolio of stocks and bond is summarized by its Sharpe ratio: expectation divided by standard deviation of return. There exists a unique portfolio of stocks with maximum Sharpe ratio, called the tangency portfolio.

The idea of beta pricing originated with the Capital Asset Pricing Model (CAPM). The CAPM assumes an equilibrium in which agents care only about mean and variance of returns, and where there are no other assets than stocks and bond. It implies that every agent holds (a scalar multiple of) the tangency portfolio. By aggregation, the market portfolio is the tangency portfolio. This result is mathematically equivalent to the exact beta pricing equation:

$$E[\tilde{r}_i] = \beta_{iM}\tau_M, \quad (3.1)$$

where $E[\cdot]$ denotes expectation, \tilde{r}_i is the return on the i^{th} stock in excess of the riskfree rate ($i = 1, \dots, N$), β_{iM} is the beta of the i^{th} stock with respect to the market, and τ_M is the expected return on the market portfolio in excess of the riskfree rate, also called risk premium. Subsequently, the idea of beta pricing was also developed outside the CAPM.

The Arbitrage Pricing Theory (APT) assumes that portfolios of stocks cannot have arbitrarily high Sharpe ratios. In other words, the maximum Sharpe ratio is finite. Since the number of stocks is finite, this assumption is trivially verified, so it must be strengthened to have economic content. The natural step is to specify *how finite* the maximum Sharpe ratio is. It means ruling out by assumption the existence of δ -arbitrage opportunities, defined as portfolios with Sharpe ratio above the predetermined level δ . This approach, called the

δ -Arbitrage Pricing Theory, is developed in the second chapter of this thesis.

Intuitively, if δ is not large then there is a tight relationship between risk and return: return only rewards those who bear risk. In this case, betas with respect to major risk factors should explain expected returns.

Let $\tilde{f}_1, \dots, \tilde{f}_K$ denote factors. We can project stock returns onto them:

$$\tilde{r}_i = \beta_{i1}\tilde{f}_1 + \dots + \beta_{iK}\tilde{f}_K + \tilde{e}_i, \quad (3.2)$$

where $\beta_{i1}, \dots, \beta_{iK}$ are factor loadings (also called betas), and \tilde{e}_i is the residual. By definition, residuals are uncorrelated with factors. Let $\bar{\lambda}$ denote the largest eigenvalue of the covariance matrix of the residuals $(\tilde{e}_i)_{i=1, \dots, N}$. If $\bar{\lambda}$ is not large, then residuals do not explain much of the risk: factors explain a lot of it.

Conversely, factors can be projected onto stock returns:

$$\tilde{f}_k = m_{k1}\tilde{r}_1 + \dots + m_{kN}\tilde{r}_N + \tilde{\eta}_k, \quad (3.3)$$

where the projection residual $\tilde{\eta}_k$ is uncorrelated with asset returns. The coefficients m_{ki} are weights of factor-mimicking portfolios. The projection residual $\tilde{\eta}_k$ has zero variance if and only if the k^{th} factor is spanned by asset returns. Let δ_F denote the maximum Sharpe measure among the portfolios that are spanned by the K factor-mimicking portfolios.

The object of interest is the approximate beta pricing equation:

$$\mathbb{E}[\tilde{r}_i] \approx \beta_{i1}\tau_1 + \dots + \beta_{iK}\tau_K, \quad (3.4)$$

where τ_1, \dots, τ_K are factor risk premia. Deviations from beta pricing are measured by the mean squared approximation error:

$$\varepsilon^2 = \min_{\tau_1, \dots, \tau_K} \frac{1}{N} \sum_{i=1}^N (\mathbb{E}[\tilde{r}_i] - \beta_{i1}\tau_1 - \dots - \beta_{iK}\tau_K)^2. \quad (3.5)$$

The smaller ε , the more accurate beta pricing. The key result is that no δ -arbitrage implies

an upper bound on deviations from beta pricing:

$$\varepsilon^2 \leq \frac{\bar{\lambda} (\delta^2 - \delta_F^2)}{N}. \quad (3.6)$$

If factors explain most of the risk in the economy ($\bar{\lambda}$ small), if approximate arbitrage opportunities are ruled out (δ not large), and if there are many stocks (N large), then betas drive most of the expected return in the economy (ε small).

3.1.2 Objective

This paper computes the bound on the right hand side of Equation (3.6) to determine how accurate beta pricing can be. Major choices must be made, for example: What is δ ? And what are the factors?

I do not wish to estimate the value of δ in this paper. I am more interested in the *mapping* between values of δ and deviations from beta pricing. Do reasonable values of δ imply that beta pricing is accurate enough to be useful? Following Section 2.3.6 of the second chapter and earlier authors, I choose *a priori* $\delta = \sqrt{2}\delta_M$, where δ_M is the Sharpe ratio of a value-weighted market index.

Two different sets of factors are investigated. The first set is a single factor equal to the return on a value-weighted index. This brings Equation (3.4) very close to the CAPM's Equation (3.1). This choice sheds some light on the accuracy of the CAPM's beta pricing equation when the CAPM's stringent equilibrium assumption is relaxed. The second set of factors contains returns on portfolios whose weights are eigenvectors corresponding to the largest eigenvalues of the covariance matrix of stock returns. This choice follows Chamberlain and Rothschild's (1983) (hereafter CR) traditional formulation of the APT. A wide range is investigated for the number of factors.

There is an interesting relationship between the two sets of factors. The market factor is highly correlated with the factor corresponding to the top covariance matrix eigenvector. Therefore a one-factor APT à la Chamberlain and Rothschild closely resembles the CAPM. As we saw, the theoretical justifications are fundamentally different. As we also saw, an important advantage of the APT formulation is that it provides a realistic bound on

deviations from beta pricing assumed away by the CAPM. But in addition, it could be said that the APT lets us bring in extra factors after the first one, in order to increase accuracy of beta pricing or, from another point of view, increase robustness to departures from the CAPM's equilibrium assumption. Therefore an important question is how much it really helps to add extra factors beyond the first one.

3.1.3 Sources of Deviation from Beta Pricing

The first source of deviation from beta pricing is obviously the riskiness of residuals, captured by $\bar{\lambda}$. The traditional APT assumes that it is negligible because residuals are infinitely less risky than factors. With a finite number of stocks, this can hardly be the case. Is the distinction between factors and residuals anywhere as extreme as the traditional APT assumes? If it is not, the bound in Equation (3.6) could be large for reasonable choices of δ .

The second source of deviation from beta pricing is risk premium estimation error. This is of particular concern if factors come from eigenvectors, in which case the sample size is limited by the number of years over which the covariance matrix of stock returns can be assumed to be stationary, usually five to twenty years. If factors are aggregate variables specified exogenously, then their risk premium is still estimated with error, but much less so because the stationarity assumption remains credible over longer horizons, possibly fifty years or more.

Assume that factors are determined by the top eigenvectors of the covariance matrix of stock returns. Call $\lambda_1, \dots, \lambda_N$ the eigenvalues of the covariance matrix. If the number of factors is K , then the maximum residual eigenvalue is $\bar{\lambda} = \lambda_{K+1}$. Let $\hat{\tau}_k$ denote the estimate of the k^{th} factor risk premium obtained from T iid observations. Neglecting estimation error on second moments, the error bound on beta pricing is:

$$\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \hat{\tau}_k \right)^2 \right] \leq \frac{\lambda_{K+1} (\delta^2 - \delta_F^2)}{N} + \frac{1}{NT} \sum_{k=1}^K \lambda_k. \quad (3.7)$$

The first term on the right hand side of Equation (3.7) determines the accuracy of the beta pricing equation, and the second term characterizes how parsimonious the set of factors

is. Accuracy increase in the number of factors K , while parsimony decreases in K . The trade-off between accuracy and parsimony can be used to find the optimal number of factors K .

3.1.4 Data

Stock returns are extracted from the Center for Research on Security Prices (CRSP) database. In order to maximize the number of stocks N , the universe consists of all the stocks traded on the New York Stock Exchange (NYSE) or the American Stock Exchange (AMEX) with less than 10% missing observations over the period considered. In order to minimize estimation error on the covariance matrix of stock returns, a high sampling frequency is chosen: daily. Nevertheless, the first and the second moments of returns are always quoted on an annual basis. To minimize estimation error on eigenvector risk premium estimates, a long period is chosen: 20 years. The period covers the first 20 years of the CRSP daily database (July 1962 to June 1982). The last twenty years might have been preferable, but they contain the Crash of 1987, an outlier which could have severely affected second moments.

The data contain 5017 daily returns on 1019 stocks. In addition, the CRSP value-weighted NYSE and AMEX index return including dividends is used as the market factor. As Roll (1977) points out, it is only a proxy, but its use is justified here since we assume that it is only *approximately* mean-variance efficient.

At the daily frequency, non-synchronous trading is an important issue. I apply Korkie's (1989) refinement of Shanken's (1987) technique to adjust covariance and beta estimates for non-synchronous trading. Cross-autocorrelation effects up to the third lag are accounted for.

3.2 Exogenous Market Factor

Some authors use the Arbitrage Pricing Theory with exogenously specified aggregates as factors. There is little consensus on the nature of these factors or even on their number, except that the return on a market proxy must be present. This is because the market return is known to explain a lot of the risk of stock returns. As mentioned above, the factor

corresponding to the top covariance matrix eigenvalue is highly correlated with the market factor.

In this section, I choose a single exogenous factor equal to the return on a value-weighted market proxy. This is especially interesting because it brings the beta pricing equations in the APT and the CAPM very close together. The spirit of this section is to find out how accurate the CAPM's beta pricing equation is when the CAPM's unrealistic equilibrium is relaxed into ruling out δ -arbitrage.

3.2.1 Factor vs. Residuals

The market factor explains 8.5% of the variance of stock returns. This is a substantial fraction, but it still leaves a lot of residual risk unexplained.

We can decompose the covariance matrix of stock returns into the part that is explained by the market factor and the part that is explained by residuals. The first one has a single nonzero eigenvalue, which is equal to 14.6 (quoted on an annual basis). The second one has $N - 1$ nonzero eigenvalues, the largest of which is equal to 3.9. Therefore factor risk easily dominates residual risk, in conformity with the intuitive properties of factors and residuals. This may be a far cry from CR's world, where the factor eigenvalue is infinitely larger than the largest residual eigenvalue, but at least it is over three times as large.

For comparison's sake, the largest covariance matrix eigenvalue is equal to 17.1 and the second largest one is equal to 1.9. This means that we could have increased the gap between factor and residuals by taking the top eigenvector instead of an exogenous market factor. However, the improvement would not have been spectacular. One advantage of the exogenous factor is that it has a clear economic interpretation.

In conclusion, this choice creates an empirical distinction between factor and residuals that satisfactorily embodies CR's theoretical distinction.

3.2.2 Residual Risk

The values for the elements entering the bound on the right hand side of Equation (3.6) are as follows: $\bar{\lambda} = 3.9$, $\delta = \sqrt{2}\delta_M$, $\delta_{F^*} = \delta_M$, $N = 1,019$. The Sharpe ratio of the

market portfolio is estimated over the whole daily CRSP database (July 1962 to December 1993): $\delta_M = 0.47$. These values correspond to an upper bound $\varepsilon \leq 2.9\%$. It means that no δ -arbitrage is compatible with deviations from CAPM beta pricing of $\pm 2.9\%$ on every expected stock return (quoted on an annual basis). This is fairly large. A 95% confidence interval for the expected stock return centered around the CAPM's prediction would be 11.4% wide.

It is instructive to compare this with the cross-sectional dispersion of CAPM predictions $\beta_{iM}\tau_M$. Their cross-sectional standard deviation is 3.1%, with the value of $\tau_M = 7.8\%$ estimated over 1926-1993 using the CRSP monthly database. Therefore the fraction of the cross-sectional dispersion of expected returns that betas do not explain is almost as large as the one that they do explain! As a consequence, ruling out δ -arbitrage opportunities is nearly consistent with a flat relationship between expected returns and betas. Unless one believes in a δ -APT with a much lower δ than the one used here, the CAPM's beta pricing equation may be too inaccurate to be of any practical use. This is due to the riskiness of CAPM residuals.

The interpretation is that, even though the market-clearing condition implies that CAPM residuals should not be priced, they are so risky that they could be priced in a way that does not create approximate arbitrage opportunities.

3.2.3 Risk Premium Estimation Error

In the most optimistic case, we can estimate the market risk premium $\hat{\tau}_M$ over the whole monthly CRSP database (January 1926 to December 1993). The standard deviation of the estimation error on $\hat{\tau}_M$ would then be quite small: 2.3%. Nevertheless, it increases the inaccuracy of CAPM beta pricing from 2.9% to 3.7%. Even when the market risk premium is estimated over a 68-year period, its estimation error manages to increase the inaccuracy of CAPM beta pricing.

3.2.4 Overall Evaluation

The APT with one exogenous factor chosen as the return on a market proxy displays a clear distinction between factors and residuals, but its beta pricing equation can be quite inaccurate. Accounting for riskiness of residuals and risk premium estimation error, the bound on the standard deviation from beta pricing is 3.7%. It may be too inaccurate for practical use. This conclusion applies to the CAPM's pricing equation too, since it is identical to the APT's in this particular case. The conclusion does depend on the value of δ for which δ -arbitrage is ruled out, but only a value much lower than the one considered here would overturn the conclusion. In my view, this is not possible without restricting investor behavior in an unrealistic way.

3.3 Eigenvectors

Chamberlain and Rothschild consider the APT with endogenously specified factors: returns on portfolios whose weights are eigenvectors corresponding to the top eigenvalues of the covariance matrix of stock returns. This choice of factors minimizes the largest residual eigenvalue $\bar{\lambda}$ under the constraint that the number of factors remains constant. However, the number of factors K is not given by CR: in their model, it is revealed by the risk structure of the stock market. I investigate a wide range of values for K , from 1 to 100.

3.3.1 Factors vs. Residuals

Figure F-10 plots the top 100 eigenvalues of the covariance matrix of stock returns. According to CR, there should be an obvious gap between the top eigenvalues (factors) and the remaining ones (residuals). According to CR too, from this point on remaining eigenvalues are negligible. These requirements on the value of K are not compatible. On the one hand, a gap is apparent only for small values of K , as Figure F-11 shows. On the other hand, the higher K goes, the more negligible residuals become, as is apparent from Figure F-12. There does not seem to be a value of K that satisfies both requirements.

Intuitively, this may mean that CR's world is not ours. Therefore the APT, which has

bite in CR's world, may not have bite in ours, for lack of a clear-cut distinction between factors and residuals. In my view, this state of affairs is responsible for the disagreements in the literature about the value of K , reviewed by Connor and Korajczyk (1992).

3.3.2 Residual Risk

Here we face an extra problem compared to Section 3.2: What is δ_F ? It is possible to estimate δ_F in the sample, but Jobson and Korkie (1980) report that it is severely biased upwards. Since the first factor is highly correlated with the market factor, we know that δ_F is of the same order as δ_M or larger. In the absence of solid evidence that it is larger, I simply take $\delta_F = \delta_M$.

With this choice, the standard deviation from beta pricing is plotted in Figure F-13. It is equal to 2.0%, 1.7% and 1.6% respectively for $K = 1, 2$ and 3. For values of K beyond 3, it decays slowly and continuously towards zero. For example, a bound of 1% can be attained by taking $K = 68$ factors. This is an order of magnitude more factors and an order of magnitude more deviation than proponents of the APT typically mention. Nevertheless, taking 68 factors is almost acceptable, considering that there are over a thousand stocks. And a formula for expected returns with standard error 1% is accurate enough for practical purposes.

Although the level of residual risk is disappointingly high, it is still low enough for the δ -APT to have an interesting implication.

3.3.3 Risk Premium Estimation Error

Using Equation (3.7), Figure F-14 plots the pricing deviations due to residual risk, risk premium estimation error, and total. Large values of K are heavily penalized by risk premium estimation error. It makes error bounds jump to 3.52%, 3.51% and 3.52% for $K = 1, 2$ and 3 respectively. For values of K beyond 3, it grows slowly and continuously. Total error is minimized by taking $K = 2$. Even the optimal choice of K yields a beta pricing equation that is not accurate enough to be of practical use.

Compared to the CAPM's 3.7%, an error of 3.52% is only a minute improvement.

Adding factors beyond the market does not really help.

3.3.4 Overall Evaluation

The APT with factors corresponding to covariance matrix eigenvectors does not display an unambiguous distinction between factors and residuals, and its beta pricing equation can be quite inaccurate. The error bound obtained here stands in sharp contrast to the back-of-the-envelope computations published by proponents of the APT. The reason is that such computations underestimate residual risk and ignore risk premium estimation error. Of course, my results could be overturned by lowering δ but, in my view, that would impose unrealistically stringent restrictions on investor behavior.

3.4 Conclusion

The link between deviations from beta pricing and approximate arbitrage opportunities in the Arbitrage Pricing Theory is weak. In practice, large deviations from beta pricing are perfectly compatible with the absence of approximate arbitrage opportunities.

Future research along the lines of this paper would examine exogenous aggregate factors other than the market factor, and investigate the impact of covariance matrix estimation error.

From a broader perspective, it is quite disappointing to find out that the trade-off between risk and return, at the present stage, does not yield useful restrictions on expected stock returns. Part of it is due to the fact that expected returns on high-risk factors, which are a key ingredient to the approach, are obviously estimated with high error.

By using the maximum residual eigenvalue, the APT obtains a weak bound that holds even in the worst-case scenario. By using the average residual eigenvalue instead, it may be possible to obtain a tighter bound that would hold in the typical-case scenario. Of course, full confidence in the bound would have to be sacrificed.

Appendix A

Spectral Theory of Large Random Matrices

This appendix gives details about the spectral theory of large-dimensional random matrices. To our knowledge, it is the first time that this theory has been mentioned in the finance literature. It bears directly on tests for the number of factors in the Arbitrage Pricing Theory (APT) based on the largest eigenvalues of the sample covariance matrix. Since this is somewhat outside the scope of the chapter, we do not provide proofs.

A.1 Mathematical Tools

A *cumulative distribution function (c.d.f.)* is a nondecreasing right-continuous function defined on the real line whose limit is zero at $-\infty$ and one at $+\infty$.

Definition 4 Let S be a symmetric matrix. Its spectral c.d.f. is the function defined by $F^S(x) = \text{proportion of eigenvalues of } S \leq x$. If the matrix S is random, so is the value of its spectral c.d.f. $F^S(x)$.

The spectral c.d.f. is in one-to-one correspondence with the system of eigenvalues. It is a convenient way to summarize the behavior of eigenvalues without invoking the joint density. The joint density would become very complicated as the number of eigenvalues grows.

Definition 5 The linear operator L transforms the c.d.f. F with support $[0, +\infty)$ into the nondecreasing function: $LF(x) = \int_{-\infty}^x t dF(t)$.

The inversion formula is: $F(x) = L^{-1}[LF](x) = LF(1) + LF(x)/x + \int_1^x LF(t) dt/t^2$ for $x > 0$, $F(0) = L^{-1}[LF](0) = \lim_{x \searrow 0} F(x)$, and $F(x) = L^{-1}[LF](x) = 0$ for $x < 0$. This linear operator is only introduced to simplify equations. Its presence can often be ignored when thinking of the problem intuitively.

Definition 6 If F is a nondecreasing function verifying $\int_{-\infty}^{+\infty} dF(t)/(1 + |t|) < \infty$, then its Stieltjes transform s_F is defined by:

$$s_F(z) = \int_{-\infty}^{+\infty} \frac{dF(t)}{t - z} \quad (\text{A.1})$$

for z on the strict upper half \mathbb{C}^+ of the complex plane. Where possible, extend s_F by continuity to real x : $s_F(x) = \lim_{z \in \mathbb{C}^+ \rightarrow x} s_F(z)$.

The inversion formula is $F(t) = \lim_{\varepsilon \searrow 0} (1/\pi) \text{Im}[\int_{-\infty}^t s_F(x + i\varepsilon) dx]$ at all points of continuity of F , where Im denotes the imaginary part of a complex number. If F is regular enough at x , e.g. twice differentiable in a neighborhood of x , then $s_F(x)$ exists and is equal to $\lim_{\varepsilon \searrow 0} \int_{|t-x| \geq \varepsilon} dF(t)/(t - x) + i\pi F'(x)$, where prime denotes the derivative (no confusion with the transposition is possible). The real and imaginary parts of s_F satisfy the Laplace equation over \mathbb{C}^+ :

$$\frac{\partial^2 \text{Re}[s_F(x + iy)]}{\partial x^2} + \frac{\partial^2 \text{Re}[s_F(x + iy)]}{\partial y^2} = 0 \quad (\text{A.2})$$

$$\frac{\partial^2 \text{Im}[s_F(x + iy)]}{\partial x^2} + \frac{\partial^2 \text{Im}[s_F(x + iy)]}{\partial y^2} = 0' \quad (\text{A.3})$$

where Re denotes the real part. For fixed $y > 0$, the function $x \mapsto (1/\pi) \text{Im}[s_F(x + iy)]$ is the convolution of the density $F'(x)$ with the Cauchy kernel $x \mapsto (y/\pi)/(x^2 + y^2)$.

Definition 7 The c.d.f.'s $(F_n)_{n \geq 1}$ converge in distribution to F if $F_n(x) \rightarrow F(x)$ at all points of continuity of F .

With these mathematical tools, we can expose the results of the spectral theory of large-dimensional random matrices that are relevant to some tests of the APT.

A.2 Asymptotic Results

Recall that $Y = U'X$ is an $N \times T$ matrix of T iid observations on a system of N uncorrelated random variables that spans the same space as the original system. Let $(y_{11}, \dots, y_{N1})'$ denote the first column of the matrix Y . y_{11}, \dots, y_{N1} are uncorrelated with variances $\lambda_1, \dots, \lambda_N$ respectively. We need to strengthen Assumption 3.

Assumption 5 $y_{11}/\sqrt{\lambda_1}, \dots, y_{N1}/\sqrt{\lambda_N}$ are iid.

We maintain Assumption 5 throughout the remainder of this appendix. The following theorem was first proven by Marčenko and Pastur (1967). It was later generalized by a number of authors. The latest and most general version is by Silverstein (1994).

Proposition 2 *Assume that the ratio N/T converges to a finite positive limit c called the concentration. Assume that the spectral c.d.f. F^Σ of the true covariance matrix Σ converges in distribution to a c.d.f. H . Then the spectral c.d.f. $F^{\tilde{\Sigma}}$ of the sample covariance matrix $\tilde{\Sigma}$ converges almost surely in distribution to a nonrandom c.d.f. G .*

The fact that the sample spectral c.d.f. $F^{\tilde{\Sigma}}$ is asymptotically nonrandom is quite remarkable. Even though $\tilde{\Sigma}$ randomly moves around its expectation Σ , its eigenvalues remain the same (in some sense). The error on sample eigenvalues is all bias and no variance. Bias comes from the fact that G is different from H .

Basic qualitative properties are established by Silverstein and Choi (1994).

Proposition 3 *G is uniquely determined by H and c . H is uniquely determined by G and c . G converges in distribution to H as c goes to zero. G has a continuous derivative, except possibly at zero. The masses $G\{0\}$ and $H\{0\}$ that G and H respectively place at zero are related by: $G\{0\} = \max(H\{0\}, 1 - 1/c)$.*

The particular shape of the distribution of the random variables X does not matter, except through the covariance matrix Σ . Under standard asymptotics, c is zero: sample and true eigenvalues coincide. Even though the distribution of true eigenvalues need not be smooth (e.g. for $\Sigma = I$ it is discontinuous at one), the distribution of sample eigenvalues must be, except possibly at zero. Intuitively, the error of sample covariance matrix eigenvectors

smoothes out sample eigenvalues. If H places some mass at zero, then G places at least the same mass at zero. Intuitively, true eigenvalues at zero do not get smoothed out because the observed variance of their corresponding eigenvectors is exactly zero in every sample. If $c > 1$, then $\tilde{\Sigma}$ is rank-deficient, therefore it can have more eigenvalues equal to zero than Σ .

The equation linking H to G is due to Marčenko and Pastur (1967):

$$\forall z \in \mathbb{C}^+ \quad s_{LH} \left(\frac{z}{1 - c s_{LG}(z)} \right) = s_{LG}(z). \quad (\text{A.4})$$

It is our contribution to introduce the linear operator L . It simplifies the equation. Equation (A.4) clearly displays how nonzero concentrations drive G and H apart. An additional advantage is that s_{LG} and s_{LH} are better behaved near zero than the Stieltjes transforms s_G and s_H used previously.

Yin (1986) derives another equation with H and G .

Proposition 4 *Assume that all the moments h_1, h_2, \dots of H exist and satisfy Carleman's condition $\sum_{k=1}^{\infty} h_{2k}^{-1/2k} = +\infty$. Then all the moments g_1, g_2, \dots of G exist and satisfy Carleman's condition. They are given by:*

$$\forall k = 1, 2, \dots \quad g_k = \sum_{w=1}^k c^{k-w} \sum \frac{k!}{n_1! n_2! \dots n_w!} h_1^{n_1} h_2^{n_2} \dots h_k^{n_k}, \quad (\text{A.5})$$

where the inner sum extends over all w -tuples of nonnegative integers (n_1, n_2, \dots, n_w) such that $\sum_{i=1}^w n_i = k - w + 1$ and $\sum_{i=1}^w i n_i = k$.

Carleman's condition ensures that a distribution is uniquely determined by its moments. It is verified by most familiar distributions whose moments exist. For the first moment, Equation (A.5) yields $g_1 = h_1$, a result that we have already seen in Theorem 2. For the second moment, $g_2 = h_2 + c h_1^2$, a result that we have already seen in Footnote 2. The second and higher moments of the sample spectral c.d.f. are larger than those of the true spectral c.d.f. The difference increases in the concentration. This means that sample eigenvalues are more dispersed than true ones. Excess dispersion increases in the concentration.

A.3 From True to Sample Eigenvalues

For $\Sigma = I$, all eigenvalues are equal to one. The true spectral c.d.f. is $H(x) = \mathbf{I}_{[0, +\infty)}(x)$, where \mathbf{I} denotes the indicator function of a set. Marčenko and Pastur solve Equation (A.4) explicitly in this important particular case. Define $a_c = (1 - \sqrt{c})^2$ and $b_c = (1 + \sqrt{c})^2$. Let $\psi_c(t) = \sqrt{(t - a_c)(b_c - t)/(2\pi ct)}$ for $a_c \leq t \leq b_c$ and $\psi_c(t) = 0$ otherwise. Then $G(x) = \int_{-\infty}^x \psi_c(t) dt$ if $0 < c \leq 1$, and $G(x) = (1 - 1/c)\mathbf{I}_{[0, +\infty)}(x) + \int_{-\infty}^x \psi_c(t) dt$ if $c > 1$. This is the formula that yields Figure F-1.

In the general case, remember that the sample spectral c.d.f. G has a continuous derivative G' , except possibly at zero. Silverstein and Choi (1994) show that for every $x \neq 0$ for which $G'(x) > 0$, $\pi c G'(x)$ is the imaginary part of the unique $z \in \mathbb{C}^+$ satisfying:

$$x = -\frac{1}{z} + c \int_{-\infty}^{+\infty} \frac{t}{1 + tz} dH(t). \quad (\text{A.6})$$

When H is discrete and its support has a finite number of points n_H , z is the root of a polynomial of degree at most $n_H + 1$. For $n_H \leq 3$, the polynomial equation can be solved in closed form, which yields an explicit formula for $G'(x)$. A Fortran routine by Wachter (1976) implements it for $n_H = 2$. Otherwise, it is straightforward to solve Equation (A.6) numerically. In particular, it is a well-posed problem.

A.4 From Sample to True Eigenvalues

The APT makes assumptions about the eigenvalues of the true covariance matrix Σ of the returns on all stocks traded in the stock market (Chamberlain and Rothschild, 1983). Some authors have tried to test these assumptions by using the eigenvalues of the sample covariance matrix $\tilde{\Sigma}$. As Brown (1989) points out and our analysis confirms, sample eigenvalues do not estimate true eigenvalues well when N is of the same order of magnitude as T , which is the usual case. In particular, the largest sample eigenvalues are upward biased estimators of the largest true eigenvalues. How can we use the spectral theory of large-dimensional matrices for such tests?

Theorem 3 states that the true spectral c.d.f. H is uniquely determined by the sample

spectral c.d.f. G and the concentration c . It is easy to obtain a smooth nonparametric estimator \hat{G} of G . Can we plug it, along with $c = N/T$, into Equation (A.4) in order to back up an estimator \hat{H} of H ?

\hat{G} can be used to estimate the complex function s_{LG} by $s_{L\hat{G}}$ over C^+ . Equation (A.4) then yields an estimator $s_{L\hat{H}}$ of the complex function s_{LH} , but not over all of C^+ : only over the domain $\hat{D} = \{z/[1 - c s_{L\hat{G}}(z)], z \in C^+\}$. This domain is included in C^+ , but excludes a portion of C^+ near the real axis. A typical domain \hat{D} is shown in Figure F-9.

From the Stieltjes transform $s_{L\hat{H}}$, we need to back up an estimate of the distribution of true eigenvalues H . Roughly, the Stieltjes inversion formula is: $\lim_{\epsilon \searrow 0} \text{Im}\{s_{L\hat{H}}(x + i\epsilon)\} = \pi x H'(x)$, where $H'(x)$ is the density of true eigenvalues. Therefore we can estimate $H'(x)$ if we know $s_{L\hat{H}}(x + i\epsilon)$ for small $\epsilon > 0$. What we need is to extend our estimator $s_{L\hat{H}}$ from the domain \hat{D} towards the real axis.

The imaginary part of $s_{L\hat{H}}$ satisfies the Laplace equation (A.2)-(A.3) over C^+ , and in particular between \hat{D} and the real line. Our goal is to solve this partial differential equation over $C^+ - \hat{D}$. The boundary of $C^+ - \hat{D}$ is divided into two pieces: the frontier with \hat{D} , where we know the value of $s_{L\hat{H}}$, and the real axis, where we want to know it. Since we do not have any information about the function on a piece of the boundary of the domain, this p.d.e. has a "free boundary."

Solving the Laplace equation with a free boundary is an ill-posed problem.

Even infinitesimal errors on the value of $s_{L\hat{H}}$ over the domain \hat{D} are amplified into large errors near the real axis. To put it in another way, there are some very different values of $s_{L\hat{H}}$ near the real axis that imply almost the same values of $s_{L\hat{H}}$ on the domain \hat{D} . Available data do not provide much guidance in choosing between them. If $s_{L\hat{H}}$ oscillated wildly over the real line, the Laplace equation would smooth it out so that we would not notice it over \hat{D} .

In practice, for high values of c , sample eigenvalues look a lot like in Figure F-1, regardless of how true eigenvalues are distributed. It is possible to back up the average and the dispersion of true eigenvalues, but not much more than that when N is of the same order of magnitude as T .

The degree of ill-posedness is determined by how far from the real line the domain \hat{D}

is. It increases in the concentration c . If c is negligible, then the domain \widehat{D} is so close to the real line that ill-posedness is negligible. In practice, c is not negligible, which is why we want to improve over sample eigenvalues in the first place.

There is, however, one reason to hope that this approach can potentially yield APT tests: the degree of ill-posedness is not uniform. It is roughly proportional to the density of sample eigenvalues. In Figure F-9, there are a lot of small eigenvalues and a few large ones. This is realistic for the stock market. We can see that the domain \widehat{D} gets closer to the real axis around large eigenvalues (large values of x). It may make it easier to estimate the density of true eigenvalues $h(x)$ when x is large. Silverstein and Combettes (1992) make a similar argument in the context of signal detection. It suggests that the problem of estimating large, isolated eigenvalues may not be ill-posed, even if the concentration is not negligible. This suggestion will be explored in future research.

In the end, it may even turn out that large, isolated true eigenvalues are actually well estimated by large, isolated sample eigenvalues. This kind of reassurance, however, cannot come from standard asymptotics. Therefore it is essential to recognize in APT tests that the number of variables N is not negligible with respect to the number of observations T . The spectral theory of large-dimensional random matrices offers one possible way to do this. Another way is proposed by Adamek (1994). He obtains very interesting results by assuming that the number of variables N goes to infinity while the number of observations T remains fixed.

Appendix B

Other Structured Estimators

This appendix discusses the optimal combination of a structured estimator $\bar{\Sigma} = [\bar{\sigma}_{ij}]_{i,j=1,\dots,N}$ and the sample covariance matrix $\tilde{\Sigma} = [\tilde{\sigma}_{ij}]_{i,j=1,\dots,N}$. Section 1.3.3 shows the importance of $\varphi_{ij} = \text{Cov}[\bar{\sigma}_{ij}, \tilde{\sigma}_{ij}]$ for $i, j = 1, \dots, N$, and $\varphi = (1/N) \sum_{i=1}^N \sum_{j=1}^N \varphi_{ij}$. This section shows how to estimate these parameters for various choices of the structured estimators $\bar{\Sigma}$.

B.1 All Variances, Respectively Covariances, Are Equal

Frost and Savarino (1986) propose a structured estimator of the covariance matrix with two free parameters: one on the diagonal, the other one off the diagonal. They obtain $\bar{\Sigma} = \hat{m}I + \hat{q}(\mathbf{1}\mathbf{1}' - I)$, where $\hat{q} = \frac{2}{N(N-1)} \sum_{i=1}^N \sum_{j=1}^{i-1} \tilde{\sigma}_{ij}$ is the average of the off-diagonal elements of the sample covariance matrix, and $\mathbf{1}$ is a conformable column vector of ones.

On the diagonal, φ_{ii} is at most of order $1/T$ for each $i = 1, \dots, N$. Off the diagonal, $\text{Var}[\tilde{\sigma}_{ij}]$ is of order $1/T$ and $\text{Var}[\bar{\sigma}_{ij}]$ is at most of order $1/(NT)$, therefore φ_{ij} is at most of order $1/(\sqrt{NT})$ for $i, j = 1, \dots, N, i \neq j$. This makes φ at most of order \sqrt{N}/T : it vanishes asymptotically. In conclusion, for this choice of prior, we recommend $\hat{\varphi} = 0$.

B.2 Diagonal Matrix

If we impose that $\bar{\Sigma}$ is diagonal, then $\varphi_{ij} = 0$ for $i, j = 1, \dots, N, i \neq j$. Since φ_{ii} is of order $1/T$ for $i = 1, \dots, N$, this makes φ at most of order $1/T$. For this choice of prior

too, we recommend $\hat{\varphi} = 0$.

B.3 All Correlation Coefficients Are Equal

We can impose that all pairs of stock returns have the same correlation coefficient. On the diagonal, $\bar{\sigma}_{ii} = \tilde{\sigma}_{ii}$, therefore $\varphi_{ii} = \text{Cov}[\bar{\sigma}_{ii}, \tilde{\sigma}_{ii}] = \text{Var}[\tilde{\sigma}_{ii}]$, which can be estimated by $\hat{\varphi}_{ii} = (1/T^2) \sum_{t=1}^T (x_{it}^2 - \tilde{\sigma}_{ii})^2$ for $i = 1, \dots, N$, as in Theorem 7.

Let $\hat{\rho} = \frac{2}{N(N-1)} \sum_{i=1}^N \sum_{j=1}^{i-1} \tilde{\sigma}_{ij} / \sqrt{\tilde{\sigma}_{ii} \tilde{\sigma}_{jj}}$ denote the average of the correlation coefficients in the sample covariance matrix. Off the diagonal, $\bar{\sigma}_{ij} = \hat{\rho} \sqrt{\tilde{\sigma}_{ii} \tilde{\sigma}_{jj}}$. $\text{Cov}[\hat{\rho}, \tilde{\sigma}_{ij}]$ is of order at most $1/(NT)$, therefore it can be neglected. $\text{Cov}[\bar{\sigma}_{ii}, \tilde{\sigma}_{ij}]$ can be estimated by $\hat{v}_{ii,ij} = (1/T^2) \sum_{t=1}^T (x_{it}^2 - \tilde{\sigma}_{ii})(x_{it}x_{jt} - \tilde{\sigma}_{ij})$. Using the delta method, $\varphi_{ij} = \text{Cov}[\bar{\sigma}_{ij}, \tilde{\sigma}_{ij}]$ can be estimated by $\hat{\varphi}_{ij} = \hat{\rho} (\hat{v}_{ii,ij} \sqrt{\tilde{\sigma}_{jj}/\tilde{\sigma}_{ii}} + \hat{v}_{jj,ij} \sqrt{\tilde{\sigma}_{ii}/\tilde{\sigma}_{jj}}) / 2$, for $i, j = 1, \dots, N, i \neq j$. These formulas yield the estimator $\hat{\varphi} = (1/N) \sum_{i=1}^N \sum_{j=1}^N \hat{\varphi}_{ij}$ that we recommend in this case.

B.4 Single Index Model

The matrix of observations is $X = [x_{it}]_{i=1, \dots, N, t=1, \dots, T}$. On the diagonal, $\bar{\sigma}_{ii} = \tilde{\sigma}_{ii}$, therefore $\varphi_{ii} = \text{Cov}[\bar{\sigma}_{ii}, \tilde{\sigma}_{ii}] = \text{Var}[\tilde{\sigma}_{ii}]$, which can be estimated by $\hat{\varphi}_{ii} = (1/T^2) \sum_{t=1}^T (x_{it}^2 - \tilde{\sigma}_{ii})^2$ for $i = 1, \dots, N$, as in Theorem 7.

Let $[x_{Mt}]_{t=1, \dots, T}$ denote returns on the market index. Let $\tilde{\sigma}_{MM} = (1/T) \sum_{t=1}^T x_{Mt}^2$, and for $i = 1, \dots, N$, let $\tilde{\sigma}_{iM} = (1/T) \sum_{t=1}^T x_{it}x_{Mt}$. Off the diagonal, $\bar{\sigma}_{ij} = \tilde{\sigma}_{iM}\tilde{\sigma}_{jM}/\tilde{\sigma}_{MM}$. $\text{Cov}[\tilde{\sigma}_{iM}, \tilde{\sigma}_{ij}]$ can be estimated by $\hat{v}_{iM,ij} = (1/T^2) \sum_{t=1}^T (x_{it}x_{Mt} - \tilde{\sigma}_{iM})(x_{it}x_{jt} - \tilde{\sigma}_{ij})$. Similarly, $\text{Cov}[\tilde{\sigma}_{MM}, \tilde{\sigma}_{ij}]$ can be estimated by $\hat{v}_{MM,ij} = (1/T^2) \sum_{t=1}^T (x_{Mt}^2 - \tilde{\sigma}_{MM})(x_{it}x_{jt} - \tilde{\sigma}_{ij})$. Using the delta method, $\varphi_{ij} = \text{Cov}[\bar{\sigma}_{ij}, \tilde{\sigma}_{ij}]$ can be estimated by $\hat{\varphi}_{ij} = \hat{v}_{iM,ij}\tilde{\sigma}_{jM}/\tilde{\sigma}_{MM} + \hat{v}_{jM,ij}\tilde{\sigma}_{iM}/\tilde{\sigma}_{MM} - \hat{v}_{MM,ij}\tilde{\sigma}_{iM}\tilde{\sigma}_{jM}/\tilde{\sigma}_{MM}^2$, for $i, j = 1, \dots, N, i \neq j$. These formulas yield the estimator $\hat{\varphi} = (1/N) \sum_{i=1}^N \sum_{j=1}^N \hat{\varphi}_{ij}$ that we recommend when $\bar{\Sigma}$ is given by the single index model.

The extension to multiple index models is tedious but straightforward.

Appendix C

Proofs of Chapter 1

We prove the theorems contained in the main body of the text. The propositions in Appendix A and the formulas in Appendix B are not proven.

C.1 Theorem 1

Recall that the matrix Y is defined as $Y = U'X$, where U is a rotation matrix containing the eigenvectors of the covariance matrix Σ . Let $[\lambda_{ij}]_{i,j=1,\dots,N}$ denote the entries of $\Lambda = U'\Sigma U$. The rotation matrix U is such that $\lambda_{ij} = 0$ when $i \neq j$ and $\lambda_{11}, \dots, \lambda_{NN}$ are the eigenvalues of the covariance matrix Σ . Please be aware that the eigenvalues of Σ are also denoted $\lambda_1, \dots, \lambda_N$ elsewhere in the text. Let $r_2^2 = E[\|\tilde{\Sigma} - \Sigma\|^2]$.

$$\begin{aligned} r_2^2 &= E \left[\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \left(\frac{1}{T} \sum_{t=1}^T y_{it} y_{jt} - \lambda_{ij} \right)^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N E \left[\left(\frac{1}{T} \sum_{t=1}^T y_{it} y_{jt} - \lambda_{ij} \right)^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \text{Var} \left[\frac{1}{T} \sum_{t=1}^T y_{it} y_{jt} \right] \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \frac{1}{T} \text{Var} [y_{i1} y_{j1}] \\ &= \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N (E [y_{i1}^2 y_{j1}^2] - E [y_{i1} y_{j1}]^2) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{NT} \sum_{i=1}^N \left(\mathbb{E} [y_{i1}^4] - 2\mathbb{E} [y_{i1}^2]^2 \right) + \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} [y_{i1}^2] \mathbb{E} [y_{j1}^2] \\
&= \frac{1}{NT} \sum_{i=1}^N \left(\mathbb{E} [y_{i1}^4] - 2\mathbb{E} [y_{i1}^2]^2 \right) + \frac{N}{T} m^2.
\end{aligned}$$

Therefore

$$\begin{aligned}
\left| \frac{N}{T} m^2 - r_2^2 \right| &\leq \frac{1}{NT} \sum_{i=1}^N \mathbb{E} [y_{i1}^4] + \frac{2}{NT} \sum_{i=1}^N \mathbb{E} [y_{i1}^2]^2 \\
&\leq \frac{3}{T} \sqrt{\frac{1}{N} \sum_{i=1}^N \mathbb{E} [y_{i1}^8]} \\
&\leq \frac{3\sqrt{B}}{T} \rightarrow 0
\end{aligned}$$

Therefore $(N/t)m^2 - r_2^2 \rightarrow 0$. \square

C.2 Theorem 2

Let $\tilde{\Sigma} = [\tilde{\sigma}_{ij}]_{i,j=1,\dots,N}$ and $\Sigma = [\sigma_{ij}]_{i,j=1,\dots,N}$. Then it can easily be shown that the following equations hold: $\mathbb{E}[(1/N) \sum_{i=1}^N \tilde{\lambda}_i] = \mathbb{E}[(1/N) \sum_{i=1}^N \tilde{\sigma}_{ii}] = (1/N) \sum_{i=1}^N \sigma_{ii} = (1/N) \sum_{i=1}^N \lambda_i$. This proves the first statement of Theorem 2.

Now let us prove the second statement. Recall that the matrix Y is defined as $Y = U'X$, where U is a rotation matrix containing the eigenvectors of the covariance matrix Σ . Let $[y_{it}]_{i=1,\dots,N, t=1,\dots,T}$ denote the entries of the matrix Y . Let $[\lambda_{ij}]_{i,j=1,\dots,N}$ denote the entries of $\Lambda = U'\Sigma U$. The rotation matrix U is such that $\lambda_{ij} = 0$ when $i \neq j$ and $\lambda_{11}, \dots, \lambda_{NN}$ are the eigenvalues of the covariance matrix Σ .

$$\begin{aligned}
\mathbb{E}[(\widehat{m} - m)^4] &= \mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T (y_{it}^2 - \lambda_{ii})\right\}^4\right] \\
&= \mathbb{E}\left[\left\{\frac{1}{T} \sum_{t=1}^T \frac{1}{N} \sum_{i=1}^N (y_{it}^2 - \lambda_{ii})\right\}^4\right] \\
&= \frac{1}{T^4} \sum_{t_1=1}^T \sum_{t_2=1}^T \sum_{t_3=1}^T \sum_{t_4=1}^T \mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_1}^2 - \lambda_{ii})\right\} \left\{\frac{1}{N} \sum_{i=1}^N (y_{it_2}^2 - \lambda_{ii})\right\} \right. \\
&\quad \left. \times \left\{\frac{1}{N} \sum_{i=1}^N (y_{it_3}^2 - \lambda_{ii})\right\} \left\{\frac{1}{N} \sum_{i=1}^N (y_{it_4}^2 - \lambda_{ii})\right\}\right] \quad (\text{C.1})
\end{aligned}$$

In the summation on the right hand side of Equation (C.1), the expectation is nonzero only if $t_1 = t_2$ or $t_1 = t_3$ or $t_1 = t_4$ or $t_2 = t_3$ or $t_2 = t_4$ or $t_3 = t_4$. Since these six conditions are symmetric we have:

$$\begin{aligned}
&\mathbb{E}[(\widehat{m} - m)^4] \\
&\leq \frac{6}{T^4} \sum_{t_1=1}^T \sum_{t_3=1}^T \sum_{t_4=1}^T \left| \mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_1}^2 - \lambda_{ii})\right\}^2 \left\{\frac{1}{N} \sum_{i=1}^N (y_{it_3}^2 - \lambda_{ii})\right\} \left\{\frac{1}{N} \sum_{i=1}^N (y_{it_4}^2 - \lambda_{ii})\right\}\right] \right| \\
&\leq \frac{6}{T^4} \sum_{t_1=1}^T \sum_{t_3=1}^T \sum_{t_4=1}^T \sqrt{\mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_1}^2 - \lambda_{ii})\right\}^4\right]} \\
&\quad \times \sqrt{\mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_3}^2 - \lambda_{ii})\right\}^2 \left\{\frac{1}{N} \sum_{i=1}^N (y_{it_4}^2 - \lambda_{ii})\right\}^2\right]} \\
&\leq \frac{6}{T^4} \sum_{t_1=1}^T \sum_{t_3=1}^T \sum_{t_4=1}^T \sqrt{\mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_1}^2 - \lambda_{ii})\right\}^4\right]} \\
&\quad \times \sqrt{\mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_3}^2 - \lambda_{ii})\right\}^4\right]} \sqrt{\mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N (y_{it_4}^2 - \lambda_{ii})\right\}^4\right]} \\
&\leq \frac{6}{T} \sqrt{\frac{1}{T} \sum_{t_1=1}^T 16 \mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N y_{it_1}^2\right\}^4\right]} \sqrt{\frac{1}{T} \sum_{t_3=1}^T 16 \mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N y_{it_3}^2\right\}^4\right]} \\
&\quad \times \sqrt{\frac{1}{T} \sum_{t_4=1}^T 16 \mathbb{E}\left[\left\{\frac{1}{N} \sum_{i=1}^N y_{it_4}^2\right\}^4\right]} \\
&\leq \frac{384B}{T} \rightarrow 0
\end{aligned}$$

where B is defined by Assumption 2. Therefore $\widehat{m} - m$ converges to zero in quartic mean, hence in quadratic mean and in probability. For future reference note that $m = (1/N) \sum_{i=1}^N E[y_{i1}^2] \leq \{(1/N) \sum_{i=1}^N E[y_{i1}^8]\}^{1/4} \leq B^{1/4}$, therefore m is bounded. \square

C.3 Theorem 3

We have $E[(1/N) \sum_{i=1}^N (\tilde{\lambda}_i - m)^2] = E[\|\tilde{\Sigma} - mI\|^2]$ and $E[(1/N) \sum_{i=1}^N (\lambda_i - m)^2] = E[\|\Sigma - mI\|^2]$. Note that $\Sigma - mI$ and $\tilde{\Sigma} - \Sigma$ are orthogonal in the sense that $E[(\Sigma - mI) \circ (\tilde{\Sigma} - \Sigma)] = (\Sigma - mI) \circ E[\tilde{\Sigma} - \Sigma] = (\Sigma - mI) \circ (\Sigma - \Sigma) = 0$. Therefore the triangle $(mI, \Sigma, \tilde{\Sigma})$ has a right angle at Σ . Then Theorem 3 follows from Pythagorus' Theorem. \square

C.4 Theorem 4

Let S denote an $N \times N$ symmetric matrix and V an $N \times N$ rotation matrix: $VV' = V'V = I$. First, note that $(1/N)\bar{r}(V'SV) = (1/N)\bar{r}(S)$. The average of the diagonal elements is invariant by rotation. Call it m . Let v_i denote the i^{th} column of V . The dispersion of the diagonal elements of $V'SV$ is $(1/N) \sum_{i=1}^N (v_i'Sv_i - m)^2$. Note that $(1/N) \sum_{i=1}^N (v_i'Sv_i - m)^2 + (1/N) \sum_{i=1}^N \sum_{j \neq i}^N (v_i'Sv_j)^2 = (1/N)\bar{r}[(V'SV - mI)^2] = (1/N)\bar{r}[(S - mI)^2]$ is invariant by rotation. Therefore the rotation V maximizes the dispersion of the diagonal elements of $V'SV$ if and only if it minimizes $(1/N) \sum_{i=1}^N \sum_{j \neq i}^N (v_i'Sv_j)^2$. This is achieved by setting $v_i'Sv_j$ to zero for all $i, j = 1, \dots, N, i \neq j$. In this case, $V'SV$ is a diagonal matrix, call it D . $V'SV = D$ is equivalent to $S = VDV'$. Since V is a rotation and D is diagonal, the column of V must contain the eigenvectors of S and the diagonal of D its eigenvalues. Therefore the dispersion of the diagonal elements of $V'SV$ is maximized when these diagonal elements are equal to the eigenvalues of S . This completes the proof of Theorem 4. \square

C.5 Theorem 5

First, we prove that the solution to Equation (1.6) is of the form $\hat{\Sigma} = \omega mI + (1 - \omega)\tilde{\Sigma}$ for some weight ω . Since $\hat{\Sigma}$ is the orthogonal projection of Σ onto the plane spanned by I and $\tilde{\Sigma}$, $(\hat{\Sigma} - \Sigma) \perp I$ where \perp denotes orthogonality. Since $\tilde{\Sigma}$ is an unbiased estimator of Σ and I is nonstochastic, $E[\tilde{\Sigma} \circ I] = \Sigma \circ I$ and $(\tilde{\Sigma} - \Sigma) \perp I$. Since mI is the orthogonal projection of Σ onto the line spanned by I , $(\Sigma - mI) \perp I$. Combining the last result with the first two yields $(\hat{\Sigma} - mI) \perp I$ and $(\tilde{\Sigma} - mI) \perp I$, therefore both $\hat{\Sigma} - mI$ and $\tilde{\Sigma} - mI$ belong to the orthogonal of I in the plane spanned by I and $\tilde{\Sigma}$, which is a subspace of dimension one. $\hat{\Sigma} - mI$ and $\tilde{\Sigma} - mI$ must be parallel, which means that $\hat{\Sigma}$ is on the line going from mI to $\tilde{\Sigma}$.

Now, we find the weight ω . The proof relies on elementary geometric relations in the triangle $(mI, \Sigma, \tilde{\Sigma})$ with right angle at Σ . $\hat{\Sigma}$ is the orthogonal projection of Σ onto the line going from mI to $\tilde{\Sigma}$. Let $d_1^2 = E[\|\hat{\Sigma} - mI\|^2]$ and $d_2^2 = E[\|\hat{\Sigma} - \tilde{\Sigma}\|^2]$. The cosine of the angle at mI can be expressed in two different ways: d_1/r_1 and r_1/d , therefore the two ratios must be equal and $d_1 = r_1^2/d$. Similarly the cosine of the angle at $\tilde{\Sigma}$ can be expressed in two different ways: d_2/r_2 and r_2/d , therefore the two ratios must be equal and $d_2 = r_2^2/d$. Note that $d_1 + d_2 = d$ as expected. These values for d_1 and d_2 yield $\hat{\Sigma} = (d_2/d)mI + (d_1/d)\tilde{\Sigma} = (r_2^2/d^2)mI + (r_1^2/d^2)\tilde{\Sigma}$.

Finally, we compute the mean squared error of $\hat{\Sigma}$. Let $r_0^2 = E[\|\hat{\Sigma} - \Sigma\|^2]$. The angle of $(\tilde{\Sigma}, mI, \Sigma)$ at mI and the angle of $(\hat{\Sigma}, \Sigma, \tilde{\Sigma})$ at Σ are equal. Equating their cosines yields $r_1/d = r_0/r_2$, therefore $r_0 = r_1r_2/d$ and $\|\hat{\Sigma} - \Sigma\|^2 = r_1^2r_2^2/d^2$. Note that Theorem 5 can also be proved by calculus alone. \square

C.6 Theorem 6

First, it is convenient to prove the following lemma.

Lemma 1 $E[\|\tilde{\Sigma}\|^2]$ is bounded.

Let $[\lambda_{ij}]_{i,j=1,\dots,N}$ denote the entries of $\Lambda = U'\Sigma U$. The rotation matrix U is such that $\lambda_{ij} = 0$ when $i \neq j$ and $\lambda_{11}, \dots, \lambda_{NN}$ are the eigenvalues of the covariance matrix Σ .

$$\begin{aligned}
E[\|\tilde{\Sigma} - \Sigma\|^2] &= E\left[\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \left(\frac{1}{T} \sum_{t=1}^T y_{it}y_{jt} - \lambda_{ij}\right)^2\right] \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N E\left[\left(\frac{1}{T} \sum_{t=1}^T y_{it}y_{jt} - \lambda_{ij}\right)^2\right] \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \text{Var}\left[\frac{1}{T} \sum_{t=1}^T y_{it}y_{jt}\right] \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \frac{1}{T} \text{Var}[y_{i1}y_{j1}] \\
&\leq \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N E[y_{i1}^2 y_{j1}^2] \\
&\leq \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N \sqrt{E[y_{i1}^4] E[y_{j1}^4]} \\
&\leq \frac{N}{T} \left(\frac{1}{N} \sum_{i=1}^N \sqrt{E[y_{i1}^4]}\right)^2 \\
&\leq \frac{N}{T} \sqrt{\frac{1}{N} \sum_{i=1}^N E[y_{i1}^8]} \\
&\leq A\sqrt{B}
\end{aligned}$$

where A and B are defined by Assumptions 1-2. Therefore $E[\|\tilde{\Sigma} - \Sigma\|^2]$ is bounded. $\|\Sigma\|^2 = (1/N) \sum_{i=1}^N E[y_{i1}^2]^2 \leq \{(1/N) \sum_{i=1}^N E[y_{i1}^8]\}^{1/2} \leq \sqrt{B}$ implies that $E[\|\tilde{\Sigma}\|^2]$ is bounded. For future reference note that it implies that d^2 , r_1^2 and r_2^2 are bounded too. \square

Now we turn to the proof of Theorem 6. We successively decompose $\hat{d}^2 - d^2$ into terms that are easier to study.

$$\hat{d}^2 - d^2 = \left\{ \|\tilde{\Sigma} - \widehat{m}I\|^2 - \|\tilde{\Sigma} - mI\|^2 \right\} + \left\{ \|\tilde{\Sigma} - mI\|^2 - E\left[\|\tilde{\Sigma} - mI\|^2\right] \right\} \quad (\text{C.2})$$

It is sufficient to show that both bracketed terms on the right hand side of Equation (C.2) converge to zero in quadratic mean. Consider the first term: $\|\tilde{\Sigma} - \widehat{m}I\|^2 - \|\tilde{\Sigma} - mI\|^2 = (\widehat{m} - m)^2$, therefore by the proof of Theorem 2 it converges to zero in quadratic mean. Now

consider the second term.

$$\|\tilde{\Sigma} - mI\|^2 = m^2 - 2m\widehat{m} + \|\tilde{\Sigma}\|^2. \quad (\text{C.3})$$

Again it is sufficient to show that the three terms on the right hand side of Equation (C.3) converge to their expectations in quadratic mean. The first term m^2 trivially does. The second term $2m\widehat{m}$ does too by the proof of Theorem 2, keeping in mind that m is bounded. Now consider the third term $\|\tilde{\Sigma}\|^2$.

$$\begin{aligned} \|\tilde{\Sigma}\|^2 &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \left(\frac{1}{T} \sum_{t=1}^T y_{it} y_{jt} \right)^2 \\ &= \frac{N}{T^2} \sum_{t=1}^T \sum_{\tau=1}^T \left(\frac{1}{N} \sum_{i=1}^N y_{it} y_{i\tau} \right)^2 \\ &= \frac{N}{T^2} \sum_{t=1}^T \left(\frac{1}{N} \sum_{i=1}^N y_{it}^2 \right)^2 + \frac{N}{T^2} \sum_{t=1}^T \sum_{\substack{\tau=1 \\ \tau \neq t}}^T \left(\frac{1}{N} \sum_{i=1}^N y_{it} y_{i\tau} \right)^2 \end{aligned} \quad (\text{C.4})$$

Again it is sufficient to show that both terms on the right hand side of Equation (C.4) converge to their expectations in quadratic mean. Consider the first term.

$$\begin{aligned} \text{Var} \left[\frac{N}{T^2} \sum_{t=1}^T \left(\frac{1}{N} \sum_{i=1}^N y_{it}^2 \right)^2 \right] &= \frac{N^2}{T^3} \text{Var} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1}^2 \right)^2 \right] \\ &\leq \frac{N^2}{T^3} \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1}^2 \right)^4 \right] \\ &\leq \left(\frac{1}{T} \right) \left(\frac{N}{T} \right)^2 \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E} [y_{i1}^8] \right) \\ &\leq \frac{A^2 B}{T} \rightarrow 0 \end{aligned}$$

Therefore the first term on the right hand side of Equation (C.4) converges to its expectation in quadratic mean.

Now consider the second term.

$$\begin{aligned} \text{Var} \left[\frac{N}{T^2} \sum_{t=1}^T \sum_{\substack{\tau=1 \\ \tau \neq t}}^T \left(\frac{1}{N} \sum_{i=1}^N y_{it} y_{i\tau} \right)^2 \right] \\ = \frac{N^2}{T^4} \sum_{t_1=1}^T \sum_{\substack{\tau_1=1 \\ \tau_1 \neq t_1}}^T \sum_{t_2=1}^T \sum_{\substack{\tau_2=1 \\ \tau_2 \neq t_2}}^T \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \quad (\text{C.5}) \end{aligned}$$

The covariances on the right hand side of Equation (C.5) only depend on $(\{t_1, \tau_1\} \cap \{t_2, \tau_2\})^\#$ the number of elements in the intersection of the set $\{t_1, \tau_1\}$ with the set $\{t_2, \tau_2\}$. This number can be zero, one or two. We study each case separately.

$$\underline{(\{t_1, \tau_1\} \cap \{t_2, \tau_2\})^\# = 0}$$

In this case $((1/N) \sum_{i=1}^N y_{it_1} y_{i\tau_1})^2$ and $((1/N) \sum_{i=1}^N y_{it_2} y_{i\tau_2})^2$ are independent., so their covariance is zero.

$$\underline{(\{t_1, \tau_1\} \cap \{t_2, \tau_2\})^\# = 1}$$

This case occurs $4t(t-1)(t-2)$ times in the summation on the right hand side of Equation (C.5). Each time we have:

$$\begin{aligned} \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \\ = \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i2} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i3} \right)^2 \right] \\ \leq \text{E} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i2} \right)^2 \left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i3} \right)^2 \right] \\ \leq \text{E} \left[\frac{1}{N^4} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N y_{i1} y_{i2} y_{j1} y_{j2} y_{k1} y_{k3} y_{l1} y_{l3} \right] \\ \leq \frac{1}{N^4} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N \text{E} [y_{i1} y_{j1} y_{k1} y_{l1}] \text{E} [y_{i2} y_{j2}] \text{E} [y_{k3} y_{l3}] \\ \leq \frac{1}{N^4} \sum_{i=1}^N \sum_{k=1}^N \text{E} [y_{i1}^2 y_{k1}^2] \text{E} [y_{i2}^2] \text{E} [y_{k3}^2] \\ \leq \frac{1}{N^4} \sum_{i=1}^N \sum_{k=1}^N \sqrt{\text{E} [y_{i1}^4]} \sqrt{\text{E} [y_{k1}^4]} \text{E} [y_{i2}^2] \text{E} [y_{k3}^2] \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{N^2} \left(\frac{1}{N} \sum_{i=1}^N \sqrt{\mathbb{E}[y_{i1}^4] \mathbb{E}[y_{i1}^2]} \right)^2 \\
&\leq \frac{1}{N^2} \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E}[y_{i1}^8] \right) \\
&\leq \frac{B}{N^2}
\end{aligned}$$

and

$$\begin{aligned}
&-\text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \\
&= -\text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i2} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i3} \right)^2 \right] \\
&\leq \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i2} \right)^2 \right] \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i3} \right)^2 \right] \\
&\leq \left(\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E}[y_{i1} y_{j1}]^2 \right)^2 \\
&\leq \frac{1}{N^2} \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E}[y_{i1}^8] \right) \\
&\leq \frac{B}{N^2}.
\end{aligned}$$

Therefore in this case the absolute value of the covariance on the right hand side of Equation (C.5) is bounded by B/N^2 .

$$\underline{(\{t_1, \tau_1\} \cap \{t_2, \tau_2\})^\# = 2}$$

This case occurs $2t(t-1)$ times in the summation on the right hand side of Equation (C.5). Each time we have:

$$\begin{aligned}
&\left| \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \right| \\
&= \left| \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i2} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{i1} y_{i2} \right)^2 \right] \right| \\
&\leq \frac{1}{N^4} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N |\text{Cov}[y_{i1} y_{i2} y_{j1} y_{j2}, y_{k1} y_{k2} y_{l1} y_{l2}]| \quad (\text{C.6})
\end{aligned}$$

Now consider the summation on the right hand side of Equation (C.6). When i, j, k, l are all pairwise distinct, Assumption 3 ensures: $E[y_{i1}y_{j1}y_{k1}y_{l1}] = E[y_{i1}y_{j1}]E[y_{k1}y_{l1}]$, which in turn implies:

$$\begin{aligned}
\text{Cov}[y_{i1}y_{i2}y_{j1}y_{j2}, y_{k1}y_{k2}y_{l1}y_{l2}] &= E[y_{i1}y_{i2}y_{j1}y_{j2}y_{k1}y_{k2}y_{l1}y_{l2}] \\
&\quad - E[y_{i1}y_{i2}y_{j1}y_{j2}]E[y_{k1}y_{k2}y_{l1}y_{l2}] \\
&= E[y_{i1}y_{j1}y_{k1}y_{l1}]^2 - E[y_{i1}y_{j1}]^2 E[y_{k1}y_{l1}] \\
&= 0.
\end{aligned}$$

Therefore the summation on the right hand side of Equation (C.6) only extends over the set $S = \{(i, j, k, l) : i, j, k, l = 1, \dots, N; \{i, j, k, l\}^\# \leq 3\}$, with the convention that $\{2, 2, 3, 4\}^\# = 3$.

$$\begin{aligned}
\left| \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \right| \\
\leq \frac{1}{N^4} \sum_{(i,j,k,l) \in S} |\text{Cov}[y_{i1}y_{i2}y_{j1}y_{j2}, y_{k1}y_{k2}y_{l1}y_{l2}]| \\
\leq \frac{1}{N^4} \sum_{(i,j,k,l) \in S} \sqrt{E[y_{i1}^2 y_{i2}^2 y_{j1}^2 y_{j2}^2] E[y_{k1}^2 y_{k2}^2 y_{l1}^2 y_{l2}^2]} \\
\leq \frac{1}{N^4} \sum_{(i,j,k,l) \in S} E[y_{i1}^2 y_{j1}^2] E[y_{k1}^2 y_{l1}^2] \\
\leq \frac{1}{N^4} \sum_{(i,j,k,l) \in S} \sqrt{E[y_{i1}^4] E[y_{j1}^4] E[y_{k1}^4] E[y_{l1}^4]}
\end{aligned}$$

The summation only extends over the quadruples (i, j, k, l) where $i = j$ or $i = k$ or $i = l$ or $j = k$ or $j = l$ or $k = l$. Since these six conditions are symmetric we have:

$$\begin{aligned}
\left| \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \right| \\
\leq \frac{6}{N^4} \sum_{i=1}^N \sum_{k=1}^N \sum_{l=1}^N \sqrt{E[y_{i1}^4] E[y_{i1}^4] E[y_{k1}^4] E[y_{l1}^4]} \\
\leq \frac{6}{N} \left(\frac{1}{N} \sum_{i=1}^N E[y_{i1}^4] \right) \left(\frac{1}{N} \sum_{i=1}^N \sqrt{E[y_{i1}^4]} \right)^2
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{6}{N} \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E} [y_{i1}^4] \right)^2 \\
&\leq \frac{6}{N} \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E} [y_{i1}^8] \right) \\
&\leq \frac{6B}{N}.
\end{aligned}$$

Having studied the three possible cases, we can now bound the summation on the right hand side of Equation (C.5):

$$\begin{aligned}
&\frac{N^2}{T^4} \sum_{t_1=1}^T \sum_{\substack{\tau_1=1 \\ \tau_1 \neq t_1}}^T \sum_{t_2=1}^T \sum_{\substack{\tau_2=1 \\ \tau_2 \neq t_2}}^T \left| \text{Cov} \left[\left(\frac{1}{N} \sum_{i=1}^N y_{it_1} y_{i\tau_1} \right)^2, \left(\frac{1}{N} \sum_{i=1}^N y_{it_2} y_{i\tau_2} \right)^2 \right] \right| \\
&\leq \frac{N^2}{T^4} \left\{ 4t(t-1)(t-2) \frac{B}{N^2} + 2t(t-1) \frac{6B}{N} \right\} \\
&\leq \frac{4B(1+3A)}{T} \rightarrow 0.
\end{aligned}$$

Backing up, the second term on the right hand side of Equation (C.4) converges to its expectation in quadratic mean. Backing up again, the third term $\|\tilde{\Sigma}\|^2$ on the right hand side of Equation (C.3) converges to its expectation in quadratic mean. Backing up more, the second bracketed term on the right hand side of Equation (C.2) converges to zero in quadratic mean. Backing up one last time, $\hat{d}^2 - d^2$ converges to zero in quadratic mean, hence in probability. For future reference note that, since $\|\tilde{\Sigma} - mI\|^2$ converges to its expectation d^2 in quadratic mean and since d^2 is bounded, $\mathbb{E}[\|\tilde{\Sigma} - mI\|^4]$ is bounded. \square

C.7 Theorem 7

Again we prove this theorem by successively decomposing $\hat{r}_2^2 - r_2^2$ into terms that are easier to study.

$$\begin{aligned}
\hat{r}_2^2 - r_2^2 &= \left\{ \frac{1}{T^2} \sum_{t=1}^T \|x_{\cdot t} x_{\cdot t}^\top - \Sigma\|^2 - \mathbb{E} [\|\tilde{\Sigma} - \Sigma\|^2] \right\} \\
&\quad + \left\{ \frac{1}{T^2} \sum_{t=1}^T \|x_{\cdot t} x_{\cdot t}^\top - \tilde{\Sigma}\|^2 - \frac{1}{T^2} \sum_{t=1}^T \|x_{\cdot t} x_{\cdot t}^\top - \Sigma\|^2 \right\} \quad (\text{C.7})
\end{aligned}$$

It is sufficient to show that both bracketed terms on the right hand side of Equation (C.7) converge to zero in quadratic mean. Consider the first term.

$$\begin{aligned}
\mathbb{E} \left[\left\| \tilde{\Sigma} - \Sigma \right\|^2 \right] &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} \left[\left(\frac{1}{T} \sum_{t=1}^T x_{it} x_{jt} - \sigma_{ij} \right)^2 \right] \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \text{Var} \left[\frac{1}{T} \sum_{t=1}^T x_{it} x_{jt} - \sigma_{ij} \right] \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \frac{1}{T^2} \sum_{t=1}^T \text{Var} [x_{it} x_{jt} - \sigma_{ij}] \\
&= \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N \text{Var} [x_{i1} x_{j1} - \sigma_{ij}] \\
&= \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} [(x_{i1} x_{j1} - \sigma_{ij})^2] \\
&= \mathbb{E} \left[\frac{1}{T} \left\| x_1 x_1^\top - \Sigma \right\|^2 \right] \\
&= \mathbb{E} \left[\frac{1}{T^2} \sum_{t=1}^T \left\| x_{\cdot t} x_{\cdot t}^\top - \Sigma \right\|^2 \right]
\end{aligned}$$

Therefore the first bracketed term on the right hand side of Equation (C.7) has expectation zero. For $t = 1, \dots, T$ let $y_{\cdot t}$ denote the $n \times 1$ vector holding the t^{th} column of the matrix Y .

$$\begin{aligned}
\text{Var} \left[\frac{1}{T^2} \sum_{t=1}^T \left\| x_{\cdot t} x_{\cdot t}^\top - \Sigma \right\|^2 \right] &= \frac{1}{T} \text{Var} \left[\frac{1}{T} \left\| x_1 x_1^\top - \Sigma \right\|^2 \right] \\
&= \frac{1}{T} \text{Var} \left[\frac{1}{T} \left\| y_{\cdot 1} y_{\cdot 1}^\top - \Lambda \right\|^2 \right] \\
&= \frac{1}{N^2 T^3} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N \text{Cov} [y_{i1} y_{j1} - \lambda_{ij}, y_{k1} y_{l1} - \lambda_{kl}] \\
&= \frac{1}{N^2 T^3} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N \text{Cov} [y_{i1} y_{j1}, y_{k1} y_{l1}] \\
&\leq \frac{1}{N^2 T^3} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N \sqrt{\mathbb{E} [y_{i1}^2 y_{j1}^2] \mathbb{E} [y_{k1}^2 y_{l1}^2]} \\
&\leq \frac{1}{N^2 T^3} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N \sqrt{\mathbb{E} [y_{i1}^4] \mathbb{E} [y_{j1}^4] \mathbb{E} [y_{k1}^4] \mathbb{E} [y_{l1}^4]} \\
&\leq \frac{N^2}{T^3} \left(\frac{1}{N} \sum_{i=1}^N \sqrt{\mathbb{E} [y_{i1}^4]} \right)^4
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{N^2}{T^3} \sqrt{\frac{1}{N} \sum_{i=1}^N \mathbb{E}[y_{i1}^8]} \\
&\leq \frac{A^2 \sqrt{B}}{T}
\end{aligned}$$

Therefore the first bracketed term on the right hand side of Equation (C.7) converges to zero in quadratic mean. For future reference note that, since $\mathbb{E}[\|\tilde{\Sigma} - \Sigma\|^2]$ is bounded, it implies that $\mathbb{E}[\{(1/T^2) \sum_{t=1}^T \|x_t x_t^\top\|^2\}^2]$ is bounded. Now consider the second term.

$$\begin{aligned}
&\mathbb{E} \left[\left\{ \frac{1}{T^2} \sum_{t=1}^T \|x_t x_t^\top - \tilde{\Sigma}\|^2 - \frac{1}{T^2} \sum_{t=1}^T \|x_t x_t^\top - \Sigma\|^2 \right\}^2 \right] \\
&\leq \mathbb{E} \left[\frac{1}{T^3} \sum_{t=1}^T \left\{ \|x_t x_t^\top - \tilde{\Sigma}\|^2 - \|x_t x_t^\top - \Sigma\|^2 \right\}^2 \right] \\
&\leq \mathbb{E} \left[\frac{1}{T^3} \sum_{t=1}^T \left\{ 2(\tilde{\Sigma} - \Sigma) \circ \left(x_t x_t^\top - \frac{\tilde{\Sigma} + \Sigma}{2} \right) \right\}^2 \right] \\
&\leq \mathbb{E} \left[\frac{4}{T^3} \sum_{t=1}^T \|\tilde{\Sigma} - \Sigma\|^2 \left\| x_t x_t^\top - \frac{\tilde{\Sigma} + \Sigma}{2} \right\|^2 \right] \\
&\leq \mathbb{E} \left[\frac{4}{T} \|\tilde{\Sigma} - \Sigma\|^2 \left(\frac{1}{T^2} \sum_{t=1}^T \left\| x_t x_t^\top - \frac{\tilde{\Sigma} + \Sigma}{2} \right\|^2 \right) \right] \\
&\leq \frac{4}{T} \sqrt{\mathbb{E}[\|\tilde{\Sigma} - \Sigma\|^4]} \sqrt{\mathbb{E} \left[\left(\frac{1}{T^2} \sum_{t=1}^T \left\| x_t x_t^\top - \frac{\tilde{\Sigma} + \Sigma}{2} \right\|^2 \right)^2 \right]} \quad (\text{C.8})
\end{aligned}$$

It is sufficient to show that the last two terms on the right hand side of Equation (C.8) are bounded. It is true for $\mathbb{E}[\|\tilde{\Sigma} - \Sigma\|^4]$ since $\mathbb{E}[\|\tilde{\Sigma} - mI\|^4]$ and $\|\Sigma - mI\|$ are bounded. Now consider the last term.

$$\frac{1}{T^2} \sum_{t=1}^T \left\| x_t x_t^\top - \frac{\tilde{\Sigma} + \Sigma}{2} \right\|^2 \leq \frac{2}{T^2} \sum_{t=1}^T \|x_t x_t^\top - \Sigma\|^2 + \frac{1}{2T} \|\tilde{\Sigma} - \Sigma\|^2$$

Since $\mathbb{E}[\{(1/T^2) \sum_{t=1}^T \|x_t x_t^\top\|^2\}^2]$ and $\mathbb{E}[\|\tilde{\Sigma} - \Sigma\|^4]$ are bounded, so is the last term on the right hand side of Equation (C.8). Backing up, the second term on the right hand side of Equation (C.7) converges to zero in quadratic mean. Backing up once more, $\hat{r}_2^2 - r_2^2$ converges to zero in quadratic mean, hence in probability. \square

C.8 Theorem 8

Follows trivially from the previous two theorems. \square

C.9 Theorem 9

C.9.1 $\|\hat{\hat{\Sigma}} - \hat{\Sigma}\|^2 \xrightarrow{P} 0$

As usual the subscript t , which should index all quantities unless otherwise specified, has been omitted to make notation lighter.

$$\begin{aligned} \|\hat{\hat{\Sigma}} - \hat{\Sigma}\| &= \left\| \frac{\hat{r}_2^2}{\hat{d}^2} (\hat{m} - m) I + \left(\frac{\hat{r}_1^2}{\hat{d}^2} - \frac{r_1^2}{d^2} \right) (\tilde{\Sigma} - mI) \right\| \\ &\leq |\hat{m} - m| + \left| \frac{(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2)}{\hat{d}^2 d^2} \right| \|\tilde{\Sigma} - mI\| \end{aligned} \quad (\text{C.9})$$

It is sufficient to prove that both terms on the right hand side of Equation (C.9) converge to zero in probability. The first term does by Theorem 2. Now consider the second term. Note that both its factors $|(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2)| / (\hat{d}^2 d^2)$ and $\|\tilde{\Sigma} - mI\|$ are bounded in probability, therefore it is sufficient to prove that either one of them converges to zero in probability. Since d^2 and r_1^2 are bounded by Lemma 1, we have: $(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2) \xrightarrow{P} 0$. Let S_1 denote the set of indices t such that

$$\left| \frac{(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2)}{\hat{d}^2 d^2} \right| \leq \sqrt{|(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2)|}.$$

If the set S_1 is infinite then $|(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2)| / (\hat{d}^2 d^2) \xrightarrow{P} 0$ as t tends to infinity inside the set S_1 , and so does the second term on the right hand side of Equation (C.9). If the complementary to the set S_1 is infinite then $\hat{d}^2 d^2 \leq |(\hat{r}_1^2 - r_1^2) d^2 - r_1^2 (\hat{d}^2 - d^2)|^{1/2} \xrightarrow{P} 0$ as t tends to infinity outside the set S_1 . By Theorem 6 it implies that $d \rightarrow 0$, therefore $\|\tilde{\Sigma} - mI\| \xrightarrow{P} 0$ as t tends to infinity outside the set S_1 , and so does the second term on the right hand side of Equation (C.9). Bringing together the results obtained for t inside and outside the set S_1 yields that the second term on the right hand side of Equation (C.9)

converges to zero in probability. Backing up, $\|\hat{\hat{\Sigma}} - \hat{\Sigma}\| \xrightarrow{P} 0$ and so does $\|\hat{\hat{\Sigma}} - \hat{\Sigma}\|^2$. \square

$$\mathbf{C.9.2} \quad \mathbb{E}[\|\hat{\hat{\Sigma}} - \Sigma\|^2] - \mathbb{E}[\|\hat{\Sigma} - \Sigma\|^2] \rightarrow 0$$

$$\begin{aligned} \mathbb{E} \left[\left\| \hat{\hat{\Sigma}} - \Sigma \right\|^2 - \left\| \hat{\Sigma} - \Sigma \right\|^2 \right] &= \mathbb{E} \left[\left\| \left(\hat{\hat{\Sigma}} - \hat{\Sigma} \right) \circ \left(\hat{\hat{\Sigma}} + \hat{\Sigma} - 2\Sigma \right) \right\|^2 \right] \\ &\leq \sqrt{\mathbb{E} \left[\left\| \hat{\hat{\Sigma}} - \hat{\Sigma} \right\|^2 \right]} \sqrt{\mathbb{E} \left[\left\| \hat{\hat{\Sigma}} + \hat{\Sigma} - 2\Sigma \right\|^2 \right]} \quad (\text{C.10}) \end{aligned}$$

It is sufficient to prove that the first term on the right hand side of Equation (C.10) converges to zero and that the second term is bounded. Consider the first term.

$$\begin{aligned} \mathbb{E} \left[\left\| \hat{\hat{\Sigma}} - \hat{\Sigma} \right\|^2 \right] &= \mathbb{E} \left[\left\| \frac{r_2^2}{d^2} (\hat{m} - m) I + \left(\frac{\hat{r}_1^2}{\hat{d}^2} - \frac{r_1^2}{d^2} \right) (\hat{\Sigma} - \hat{m} I) \right\|^2 \right] \\ &= \mathbb{E} \left[\frac{r_2^4}{d^4} (\hat{m} - m)^2 \right] + \mathbb{E} \left[\left(\frac{\hat{r}_1^2}{\hat{d}^2} - \frac{r_1^2}{d^2} \right)^2 \left\| \hat{\Sigma} - \hat{m} I \right\|^2 \right] \\ &\leq \mathbb{E} \left[(\hat{m} - m)^2 \right] + \mathbb{E} \left[\frac{(\hat{r}_1^2 \hat{d}^2 - r_1^2 \hat{d}^2)^2}{\hat{d}^2 d^4} \right] \quad (\text{C.11}) \end{aligned}$$

It is sufficient to show that both terms on the right hand side of Equation (C.11) converges to zero. The first term does by the proof of Theorem 2. Now consider the second term. Since $\hat{r}_1^2 \leq \hat{d}^2$ and $\hat{r}_1^2 \leq \hat{d}^2$, note for future reference that $(\hat{r}_1^2 \hat{d}^2 - r_1^2 \hat{d}^2)^2 / (\hat{d}^2 d^4) \leq 4\hat{d}^2$. Fix $\varepsilon > 0$. Let S_3 denote the set of indices t such that $d^2 \leq \varepsilon/8$. Since $\hat{d}^2 - d^2 \rightarrow 0$ in quadratic mean, $\exists T_1 \quad \forall T \quad T \geq T_1 \Rightarrow \mathbb{E}[|\hat{d}^2 - d^2|] \leq \varepsilon/8$. We have:

$$\begin{aligned} \forall T \quad T \in S_3, t \geq T_1 \Rightarrow \mathbb{E} \left[\frac{(\hat{r}_1^2 \hat{d}^2 - r_1^2 \hat{d}^2)^2}{\hat{d}^2 d^4} \right] &\leq 4\mathbb{E}[\hat{d}^2] \\ &\leq 4\mathbb{E}[|\hat{d}^2 - d^2|] + 4d^2 \\ &\leq 4\frac{\varepsilon}{8} + 4\frac{\varepsilon}{8} \\ &\leq \varepsilon. \quad (\text{C.12}) \end{aligned}$$

Since $\hat{r}_1^2 - r_1^2$ and $\hat{d}^2 - d^2$ converge to zero in quadratic mean and since r_1^2 and d^2 are bounded, $\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2 = (\hat{r}_1^2 - r_1^2) \hat{d}^2 - r_1^2 (\hat{d}^2 - d^2) \rightarrow 0$ in quadratic mean, therefore $\exists T_2 \quad \forall T \quad T \geq T_2 \Rightarrow E[(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2] \leq \varepsilon^4/1024$. Denote $\Pr(\cdot)$ the probability of an event. We have: $\forall T \quad T \notin S_3, T \geq T_2 \Rightarrow$

$$\begin{aligned}
E \left[\frac{(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2}{\hat{d}^2 d^4} \right] &= E \left[\frac{(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2}{\hat{d}^2 d^4} \middle| \hat{d}^2 \leq \frac{\varepsilon}{8} \right] \Pr \left(\hat{d}^2 \leq \frac{\varepsilon}{8} \right) \\
&\quad + E \left[\frac{(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2}{\hat{d}^2 d^4} \middle| \hat{d}^2 > \frac{\varepsilon}{8} \right] \Pr \left(\hat{d}^2 > \frac{\varepsilon}{8} \right) \\
&\leq E \left[4 \hat{d}^2 \middle| \hat{d}^2 \leq \frac{\varepsilon}{8} \right] \Pr \left(\hat{d}^2 \leq \frac{\varepsilon}{8} \right) \\
&\quad + \frac{8}{\varepsilon d^4} E \left[(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2 \middle| \hat{d}^2 > \frac{\varepsilon}{8} \right] \Pr \left(\hat{d}^2 > \frac{\varepsilon}{8} \right) \\
&\leq 4 \frac{\varepsilon}{8} + \frac{512}{\varepsilon^3} E \left[(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2 \right] \\
&\leq \frac{\varepsilon}{2} + \frac{512}{\varepsilon^3} \frac{\varepsilon^4}{1024} \\
&\leq \varepsilon.
\end{aligned} \tag{C.13}$$

Bringing together the results from Equations (C.12)-(C.13) yields:

$$\forall T \quad T \geq \max(T_1, T_2) \Rightarrow E \left[\frac{(\hat{r}_1^2 \hat{d}^2 - r_1^2 d^2)^2}{\hat{d}^2 d^4} \right] \leq \varepsilon,$$

therefore the second term on the right hand side of Equation (C.11) converges to zero. Backing up, the first term on the right hand side of Equation (C.10) converges to zero. Since $E[\|\hat{\Sigma} - \Sigma\|^2]$ is bounded, it implies that the second term on the right hand side of Equation (C.10) is bounded too. Backing up once more yields $E[\|\hat{\Sigma} - \Sigma\|^2] - E[\|\hat{\Sigma} - \Sigma\|^2] \rightarrow 0$. \square

C.9.3 $(\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2) - (r_1^2 r_2^2 / d^2) \xrightarrow{P} 0$

$$\frac{\hat{r}_1^2 \hat{r}_2^2}{\hat{d}^2} - \frac{r_1^2 r_2^2}{d^2} = \frac{(\hat{r}_1^2 \hat{r}_2^2 - r_1^2 r_2^2) d^2 - r_1^2 r_2^2 (\hat{d}^2 - d^2)}{\hat{d}^2 d^2}$$

By Theorems 6-8 and Lemma 1 the numerator on the right hand side converges to zero in probability. Let S_2 denote the set of indices t such that

$$\left| \frac{(\hat{r}_1^2 \hat{r}_2^2 - r_1^2 r_2^2) d^2 - r_1^2 r_2^2 (\hat{d}^2 - d^2)}{\hat{d}^2 d^2} \right| \leq \sqrt{|(\hat{r}_1^2 \hat{r}_2^2 - r_1^2 r_2^2) d^2 - r_1^2 r_2^2 (\hat{d}^2 - d^2)|}.$$

If the set S_2 is infinite then $(\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2) - (r_1^2 r_2^2 / d^2) \xrightarrow{P} 0$ as t tends to infinity inside the set S_2 . If the complementary to the set S_2 is infinite then $\hat{d}^2 d^2 \leq |(\hat{r}_1^2 \hat{r}_2^2 - r_1^2 r_2^2) d^2 - r_1^2 r_2^2 (\hat{d}^2 - d^2)|^{1/2} \xrightarrow{P} 0$ as t tends to infinity outside the set S_2 . By Theorem 6 it implies that $d^2 \rightarrow 0$, therefore $(r_1^2 r_2^2 / d^2) \xrightarrow{P} 0$ as t tends to infinity outside the set S_2 , and so does $(\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2)$. Bringing together the results obtained for t inside and outside the set S_2 yields $(\hat{r}_1^2 \hat{r}_2^2 / \hat{d}^2) - (r_1^2 r_2^2 / d^2) \xrightarrow{P} 0$. \square

C.10 Theorem 10

This is similar to the proof of Theorem 5. $\hat{\Sigma}$ is the orthogonal projection of Σ on the line between $\bar{\Sigma}$ and $\tilde{\Sigma}$. Let $d_1^2 = E[||\hat{\Sigma} - mI||]$, $d_2^2 = E[||\hat{\Sigma} - \tilde{\Sigma}||]$ and $r_0^2 = E[||\hat{\Sigma} - \Sigma||^2]$. The orthogonality condition $(\bar{\Sigma} - \hat{\Sigma}) \perp (\Sigma - \hat{\Sigma})$ implies $d_1^2 + r_0^2 = r_1^2$. Also, the orthogonality condition $(\tilde{\Sigma} - \hat{\Sigma}) \perp (\Sigma - \hat{\Sigma})$ implies $d_2^2 + r_0^2 = r_2^2$. Subtracting one equation from the other yields $d_1^2 - d_2^2 = r_1^2 - r_2^2$. Since $\bar{\Sigma}$, $\hat{\Sigma}$ and $\tilde{\Sigma}$ are aligned, we have $d_1 + d_2 = d$, which implies $d_1^2 - d_2^2 = d_1^2 + (d - d_1)^2 = 2d_1 d - d^2$. Therefore $2d_1 d - d^2 = r_1^2 - r_2^2$, i.e. $d_1 = (r_1^2 + d^2 - r_2^2) / 2d$. By symmetry, $d_2 = (r_2^2 + d^2 - r_1^2) / 2d$. Note that $d_1 + d_2 = d$ as expected. These values for d_1 and d_2 yield $\hat{\Sigma} = (d_2/d)\bar{\Sigma} + (d_1/d)\tilde{\Sigma} = [(r_2^2 + d^2 - r_1^2) / (2d^2)]\bar{\Sigma} + [(r_1^2 + d^2 - r_2^2) / (2d^2)]\tilde{\Sigma}$.

Appendix D

Proofs of Chapter 2

The notation common to the proofs is as follows.

The $N \times 1$ random vector $\tilde{\mathbf{r}} = (\tilde{r}_i)_{i=1,\dots,N}$ contains stock returns. The $N \times 1$ vector $\boldsymbol{\mu} = (\mu_i)_{i=1,\dots,N} = (E[\tilde{r}_i])_{i=1,\dots,N}$ contains expected stock returns. The $N \times N$ matrix $\boldsymbol{\Sigma} = (\sigma_{ij})_{i,j=1,\dots,N} = (\text{Cov}[\tilde{r}_i, \tilde{r}_j])_{i,j=1,\dots,N}$ contains variances and covariances of stock returns.

The $K \times 1$ random vector $\tilde{\mathbf{f}} = (\tilde{f}_k)_{k=1,\dots,K}$ contains factors. The $N \times K$ matrix $\mathbf{B} = (\beta_{ik})_{\substack{i=1,\dots,N \\ k=1,\dots,K}}$ contains factor loadings. $\boldsymbol{\Psi} = \text{Var}[\tilde{\mathbf{f}}]$ is the $K \times K$ variance-covariance matrix of factors. The $K \times 1$ vector $\boldsymbol{\tau} = (\tau_k)_{k=1,\dots,K}$ contains factor risk premia. The $N \times N$ matrix $\boldsymbol{\Omega} = (\omega_{ij})_{i,j=1,\dots,N} = (\text{Cov}[\tilde{e}_i, \tilde{e}_j])_{i,j=1,\dots,N}$ contains variances and covariances of residuals.

The $N \times 1$ vector $\mathbf{w} = (w_i)_{i=1,\dots,N}$ contains portfolio weights.

D.1 Theorem 11

See Chamberlain and Rothschild (1983). This result can also be obtained in a less sophisticated way by letting N go to infinity in Theorem 13 while holding δ and $\bar{\lambda}$ bounded.

□

D.2 Theorem 12

The $N \times 1$ vector \mathbf{e} contains the endowment of marketable assets. Let \tilde{h} denote the future value of all non-marketable assets. The representative agent has risk aversion A . Her utility maximization problem is equivalent to: $\max_{\mathbf{w}} \mathbf{w}'\boldsymbol{\mu} - \frac{1}{2}A(\mathbf{w}'\boldsymbol{\Sigma}\mathbf{w} + 2\mathbf{w}'\text{Cov}[\tilde{\mathbf{r}}, \tilde{h}])$. Let $\mathbf{h} = \boldsymbol{\Sigma}^{-1}\text{Cov}[\tilde{\mathbf{r}}, \tilde{h}]$ denote the portfolio whose future value most closely mimics the future value of non-marketable assets. The solution to the representative agent's utility maximization problem is: $\mathbf{w} = -\mathbf{h} + \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}/A$. The agent first hedges her exposure to non-marketable risk by shorting portfolio \mathbf{h} , and then adds a mean-variance efficient position. The market clearing condition is $\mathbf{w} = \mathbf{e}$, therefore in equilibrium: $\boldsymbol{\mu} = A\boldsymbol{\Sigma}(\mathbf{e} + \mathbf{h})$.

The maximum squared Sharpe measure in the market is given by: $\bar{\delta}^2 = \boldsymbol{\mu}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} = A^2(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e} + 2\mathbf{e}'\boldsymbol{\Sigma}\mathbf{h} + \mathbf{h}'\boldsymbol{\Sigma}\mathbf{h})$. The squared Sharpe measure of the market portfolio is: $\delta_M^2 = (\boldsymbol{\mu}'\mathbf{e})^2/(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e}) = A^2(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e} + \mathbf{e}'\boldsymbol{\Sigma}\mathbf{h})^2/(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e})$. Therefore we have:

$$\frac{\bar{\delta}^2}{\delta_M^2} - 1 = \frac{(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e})(\mathbf{h}'\boldsymbol{\Sigma}\mathbf{h}) - (\mathbf{e}'\boldsymbol{\Sigma}\mathbf{h})^2}{(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e} + \mathbf{e}'\boldsymbol{\Sigma}\mathbf{h})^2}. \quad (\text{D.1})$$

The numerator on the right hand side of Equation (D.1) is no greater than $(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e})(\mathbf{h}'\boldsymbol{\Sigma}\mathbf{h})$. In Theorem 12, it is assumed that the covariance between marketable and non-marketable assets is non-negative, which ensures that $\mathbf{e}'\boldsymbol{\Sigma}\mathbf{h} \geq 0$. Therefore the denominator on the right hand side of Equation (D.1) is at least as great as $(\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e})^2$. It implies that:

$$\frac{\bar{\delta}^2}{\delta_M^2} - 1 \leq \frac{\mathbf{h}'\boldsymbol{\Sigma}\mathbf{h}}{\mathbf{e}'\boldsymbol{\Sigma}\mathbf{e}}. \quad (\text{D.2})$$

The variance of the value of non-marketable assets σ_{NM}^2 is even higher than the variance of its projection onto the market $\mathbf{h}'\boldsymbol{\Sigma}\mathbf{h}$. Noting that $\sigma_M^2 = \mathbf{e}'\boldsymbol{\Sigma}\mathbf{e}$ completes the proof of Theorem 12. \square

D.3 Theorem 13

This proof relies heavily on the orthogonal projection \mathbf{P} onto the residual space. Formally: $\mathbf{P} = \mathbf{I}_N - \mathbf{B}(\mathbf{B}'\mathbf{B})^{-1}\mathbf{B}'$, where \mathbf{I}_N denotes the $N \times N$ identity matrix. Recall the properties

of orthogonal projections: $P^2 = P = P'$.

Since residuals are uncorrelated with factors, we have: $\Sigma = B\Psi B' + \Omega$. It implies that $P'\Sigma P = P'(B\Psi B' + \Omega)P = P'\Omega P$.

Risk premia are determined by: $\min_{\tau} (\mu - B\tau)'(\mu - B\tau)/N$. The solution to this minimization problem is: $\tau = (B'B)^{-1}B'\mu$. The mean squared error on beta pricing is: $\varepsilon^2 = (\mu - B\tau)'(\mu - B\tau)/N = \mu'P\mu/N$.

Consider the portfolio with weights $P\mu$. Its expected return is $(P\mu)'\mu = \mu'P\mu$. Its return variance is $(P\mu)'\Sigma(P\mu) = \mu'P'\Sigma P\mu = \mu'P'\Omega P\mu = (P\mu)'\Omega(P\mu) \leq \bar{\lambda} (P\mu)'(P\mu) = \bar{\lambda} \mu'P\mu$, where the inequality follows from the properties of the largest eigenvalue of the residual covariance matrix. Therefore portfolio $P\mu$'s squared Sharpe measure exceeds $(\mu'P\mu)^2/(\bar{\lambda} \mu'P\mu) = \mu'P\mu/\bar{\lambda}$. So the maximum squared Sharpe measure in the economy $\bar{\delta}^2$ must also exceed $\mu'P\mu/\bar{\lambda}$. As it turns out, this statement can be refined.

δ_F is the maximum Sharpe measure in the projection of the factor space onto asset returns. It means that there exists a portfolio whose return is uncorrelated with portfolio $P\mu$'s and whose Sharpe measure is δ_F . It is always possible to form a linear combination of two uncorrelated assets so that the squared Sharpe measure of the combination equals the sum of the squared Sharpe measures of the two assets. In this case, it is possible to form a portfolio whose squared Sharpe measure is $\mu'P\mu/\bar{\lambda} + \delta_F^2$. As a consequence, the maximum squared Sharpe measure in the economy $\bar{\delta}^2$ must be at least as high as $\mu'P\mu/\bar{\lambda} + \delta_F^2$.

It follows that $\varepsilon^2 = \mu'P\mu/N \leq \bar{\lambda} (\bar{\delta}^2 - \delta_F^2)/N$. Assumption 4 completes the proof of Theorem 13 by providing the inequality $\bar{\delta}^2 \leq \delta^2$. \square

D.4 Theorem 14

The intuition is that the regression of stock returns on factors has a better fit if factors are spanned by stock returns, all other things being equal.

$M = (m_{ki})_{\substack{k=1,\dots,K \\ i=1,\dots,N}}$ is the $K \times N$ matrix of weights of factor-mimicking portfolios. The $K \times 1$ random vector $\tilde{\eta} = (\tilde{\eta}_k)_{k=1,\dots,K}$ contains the residuals of the projection of factors onto returns. $\Theta = \text{Var}[\tilde{\eta}]$ is its $K \times K$ variance-covariance matrix. We have: $\Psi = M\Sigma M' + \Theta$.

$\tilde{\mathbf{f}}_* = \mathbf{M}\tilde{\mathbf{r}}$ is the $K \times 1$ random vector of returns on factor-mimicking portfolios. $\Psi_* = \text{Var}[\tilde{\mathbf{f}}_*]$ is the variance-covariance matrix of $\tilde{\mathbf{f}}_*$. We have: $\Psi_* = \mathbf{M}\Sigma\mathbf{M}'$ and $\Psi = \Psi_* + \Theta$, hence $\Psi \succ \Psi_*$, where the symbol \succ represents the ordering between symmetric matrices. A useful implication is that $\Psi\Psi_*^{-1}\Psi - \Psi \succ \mathbf{0}_K$, where $\mathbf{0}_K$ is the $K \times K$ null matrix.

The coefficients of the regression of stock returns on $\tilde{\mathbf{f}}_*$ are:

$$\mathbf{B}_* = \text{Cov}[\tilde{\mathbf{r}}, \tilde{\mathbf{f}}_*'] \text{Var}[\tilde{\mathbf{f}}_*]^{-1} \quad (\text{D.3})$$

$$= \text{Cov}[\tilde{\mathbf{r}}, \tilde{\mathbf{f}}_*'] \text{Var}[\tilde{\mathbf{f}}_*]^{-1} \quad (\text{D.4})$$

$$= \text{Cov}[\tilde{\mathbf{r}}, \tilde{\mathbf{f}}_*'] \Psi_*^{-1} \quad (\text{D.5})$$

$$= \text{Cov}[\tilde{\mathbf{r}}, \tilde{\mathbf{f}}_*'] \Psi^{-1} \Psi \Psi_*^{-1} \quad (\text{D.6})$$

$$= \mathbf{B} \Psi \Psi_*^{-1} \quad (\text{D.7})$$

Let Ω_* denote the covariance matrix of the residuals of the projection of stock returns on $\tilde{\mathbf{f}}_*$. We have: $\Sigma = \mathbf{B}\Psi\mathbf{B}' + \Omega = \mathbf{B}_*\Psi_*\mathbf{B}_*' + \Omega_*$. Therefore:

$$\Omega - \Omega_* = \mathbf{B}_*\Psi_*\mathbf{B}_*' - \mathbf{B}\Psi\mathbf{B}' \quad (\text{D.8})$$

$$= (\mathbf{B}\Psi\Psi_*^{-1})\Psi_*(\mathbf{B}\Psi\Psi_*^{-1})' - \mathbf{B}\Psi\mathbf{B}' \quad (\text{D.9})$$

$$= \mathbf{B}\Psi\Psi_*^{-1}\Psi\mathbf{B}' - \mathbf{B}\Psi\mathbf{B}' \quad (\text{D.10})$$

$$= \mathbf{B}(\Psi\Psi_*^{-1}\Psi - \Psi)\mathbf{B}' \quad (\text{D.11})$$

$$\succ \mathbf{0}_N. \quad (\text{D.12})$$

Therefore the largest eigenvalue of Ω exceeds the largest eigenvalue of Ω_* . This completes the proof of Theorem 14. \square

D.5 Theorem 15

This is a well-known result from matrix algebra. \square

D.6 Theorem 16

Decompose the covariance matrix of stock returns Σ into eigenvalues and eigenvectors: $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$. The diagonal elements of the $N \times N$ matrix $\mathbf{\Lambda}$ are the eigenvalues $\lambda_1, \dots, \lambda_N$ of Σ , and the off-diagonal elements are equal to zero. The eigenvalues are sorted in descending order. The column vectors of the $N \times N$ orthogonal matrix \mathbf{U} are the eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_N$ of Σ .

Let \mathbf{U}_k denote the $N \times K$ matrix containing the first K columns of \mathbf{U} . The K factors are the returns on the portfolios whose weights are the column vectors of \mathbf{U}_k : $\tilde{\mathbf{f}} = \mathbf{U}_k' \tilde{\mathbf{r}}$. The matrix of betas is: $\mathbf{B} = \Sigma \mathbf{U}_k (\mathbf{U}_k' \Sigma \mathbf{U}_k)^{-1} = \mathbf{U}_k$.

The $K \times 1$ random vector $\hat{\boldsymbol{\tau}} = (\hat{\tau}_k)_{k=1, \dots, K}$ contains estimates of risk premia based on T iid observations. The variance-covariance matrix of estimated risk premia is: $\text{Var}[\hat{\boldsymbol{\tau}}] = \text{Var}[\tilde{\mathbf{f}}]/T = \mathbf{U}_k' \Sigma \mathbf{U}_k / T$. It is the diagonal matrix containing the top K eigenvalues of Σ divided by T .

The beta pricing equation with estimated risk premia is:

$$\mu_i \approx \sum_{k=1}^K \beta_{ik} \hat{\tau}_k \quad (\text{D.13})$$

The expected sum of squared deviations from Equation (D.13) is:

$$\mathbb{E} \left[\sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \hat{\tau}_k \right)^2 \right] = \sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \tau_k \right)^2 + \sum_{i=1}^N \text{Var} \left[\sum_{k=1}^K \beta_{ik} \hat{\tau}_k \right] \quad (\text{D.14})$$

$$= \sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \tau_k \right)^2 + \sum_{i=1}^N \sum_{k=1}^K \beta_{ik}^2 \text{Var}[\hat{\tau}_k] \quad (\text{D.15})$$

$$= \sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \tau_k \right)^2 + \sum_{k=1}^K \text{Var}[\hat{\tau}_k] \quad (\text{D.16})$$

$$= \sum_{i=1}^N \left(\mu_i - \sum_{k=1}^K \beta_{ik} \tau_k \right)^2 + \sum_{k=1}^K \frac{\lambda_k}{T}. \quad (\text{D.17})$$

This completes the proof of Theorem 16. \square

Appendix E

Tables

	Structured	Shrinkage	T-Statistic
$\widehat{m}I$	20.3	10.9	7.01
B.1	20.3	10.6	7.20
B.2	16.0	9.6	8.33
B.3	13.8	9.6	6.37
B.4	11.5	9.3	4.94

Table E.1: Comparison of the Ex-Post Standard Deviations of Ex-Ante Minimum Variance Portfolios.

Standard deviations are quoted in percents on an annual basis. The portfolios are obtained using a structured estimator of the covariance matrix, or its associated shrinkage estimator. The t-statistic tests the null hypothesis that a given structured estimator and its associated shrinkage estimator yield ex-ante minimum variance portfolios with the same ex-post variance of returns. This hypothesis is rejected in all five cases. Shrinkage helps portfolio selection minimize variance.

	Structured With Hindsight	Shrinkage Without Hindsight	T-Statistic
$\widehat{m}I$	11.8	10.9	1.88
B.1	11.8	10.6	3.80
B.2	11.0	9.6	4.38
B.3	12.4	9.6	5.53
B.4	11.0	9.3	5.65

Table E.2: Comparison of the Ex-Post Standard Deviations of Minimum Variance Portfolios.

Standard deviations are quoted in percents on an annual basis. The portfolios are obtained using a structured estimator of the covariance matrix, or its associated shrinkage estimator. For structured estimators, the minimum variance portfolio is chosen ex-post among linear combinations of three portfolios that span the ex-ante mean-variance efficient set, assuming that returns are driven by beta and size only. For shrinkage estimators, the minimum variance portfolio is chosen ex-ante, without the benefit of hindsight. This makes it harder to help portfolio selection minimize variance. The t-statistic tests the null hypothesis that a given structured estimator and its associated shrinkage estimator yield minimum variance portfolios with the same ex-post variance of returns. All reject the null. The t-statistic of 1.88 is significant at the 5% level against the one-sided alternative that shrinkage helps minimize variance.

	Plain Regression	Excluding January	Including Size	1963-1992
Slope	2.33	-0.77	0.33	1.88
Standard Error	(2.27)	(2.31)	(1.93)	(3.15)
T-Statistic	1.03	-0.33	0.17	0.60

Table E.3: Predictive OLS Cross-Sectional Regression of Returns on Betas over 1936-1992. Data come from the Center for Research in Security Prices (CRSP) database. Slope estimates are quoted in percents on an annual basis. Returns are in excess of the riskfree rate. The universe for a given year includes all common stocks traded on the NYSE and (after 1963) AMEX, with all valid monthly returns over the past 10 years and valid market capitalization. Returns are buy-and-hold, with annual rebalancing.

	Plain Regression	Excluding January	Including Size	1963-1992
Slope	3.51	3.08	2.57	3.08
Standard Error	(1.84)	(1.90)	(1.78)	(2.66)
T-Statistic	1.91	1.62	1.44	1.16

Table E.4: Predictive GLS Cross-Sectional Regression of Returns on Betas over 1936-1992. Data come from the Center for Research in Security Prices (CRSP) database. Slope estimates are quoted in percents on an annual basis. Returns are in excess of the riskfree rate. The universe for a given year includes all common stocks traded on the NYSE and (after 1963) AMEX, with all valid monthly returns over the past 10 years and valid market capitalization. Returns are buy-and-hold, with annual rebalancing. The covariance matrix estimate required for GLS is obtained from the asymptotic shrinkage estimator associated with the structured estimator from Appendix B.4 (single index model).

	Plain Regression	Excluding January	Including Size	1963-1992
Section 1.3.2	2.58	2.23	1.14	2.35
	(1.82)	(1.88)	(1.73)	(2.56)
	1.42	1.19	0.66	0.92
Appendix B.1	2.53	2.06	1.14	2.22
	(1.81)	(1.86)	(1.71)	(2.55)
	1.40	1.11	0.66	0.87
Appendix B.2	3.61	3.44	2.65	3.01
	(1.82)	(1.88)	(1.80)	(2.63)
	1.98	1.83	1.47	1.14
Appendix B.3	3.39	4.56	3.42	3.85
	(1.77)	(1.80)	(1.71)	(2.45)
	1.92	2.53	2.00	1.57

Table E.5: Predictive GLS Cross-Sectional Regression of Returns on Betas over 1936-1992. Data come from the Center for Research in Security Prices (CRSP) database. In each cell, the first number is the slope estimates are quoted in percents on an annual basis; the second number (in parenthesis) is the standard error on this number; and the third number is the t-statistic obtained by dividing the above two numbers. Returns are in excess of the riskfree rate. The universe for a given year includes all common stocks traded on the NYSE and (after 1963) AMEX, with all valid monthly returns over the past 10 years and valid market capitalization. Returns are buy-and-hold, with annual rebalancing. The covariance matrix estimates required for GLS is obtained from the asymptotic shrinkage estimator associated with the structured estimator from Section 1.3.2, and Appendices B.1, B.2 and B.3 respectively.

Appendix F

Figures

Figures start on the next page.

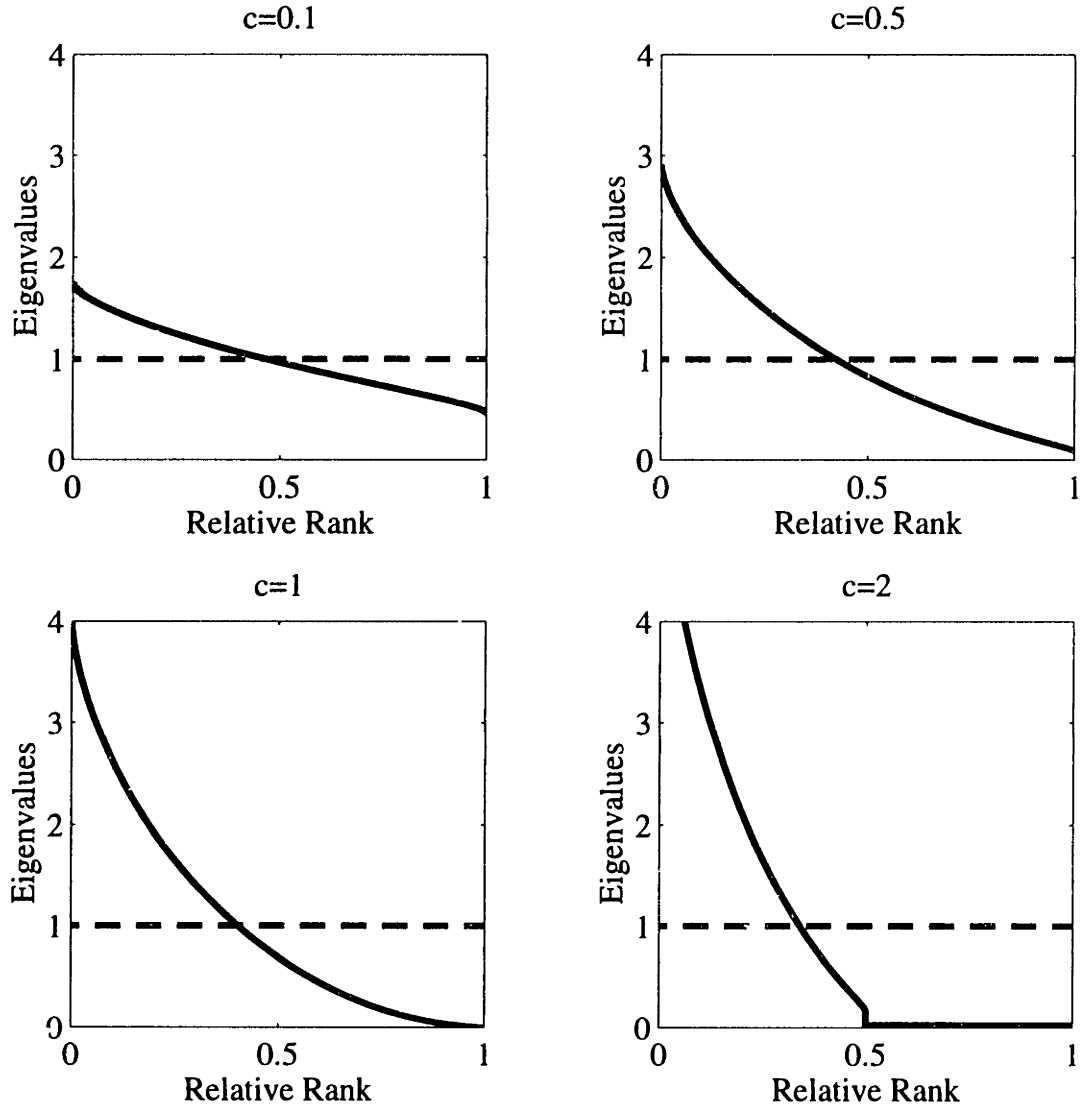


Figure F-1: Sample vs. True Eigenvalues.

The solid line represents the distribution of the eigenvalues of the sample covariance matrix. Eigenvalues are sorted in descending order, then plotted against their relative rank, defined as the ratio of the rank to the total number of eigenvalues N . When N changes, the relative rank remains between zero (largest eigenvalues) and one (smallest). We assume that the true covariance matrix is the identity, i.e. true eigenvalues are equal to one. The distribution of true eigenvalues is plotted as the dashed horizontal line. Distributions are obtained in the limit as the number of observations T and the number of variables N both go to infinity, with their ratio N/T converging to a finite positive limit c called the concentration. The four plots correspond to different concentrations. The smallest eigenvalues of the sample covariance matrix are severely biased downwards and the largest ones upwards. Bias increases in the concentration.

Geometric Interpretation of Theorem 5

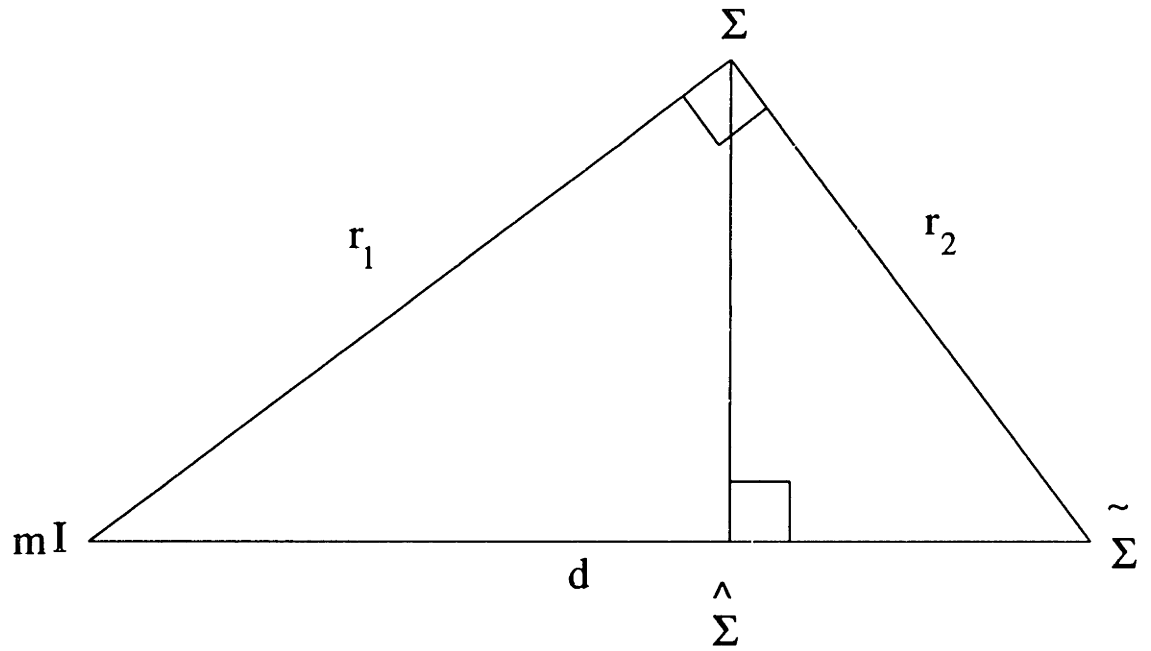


Figure F-2: Geometric Interpretation of Theorem 5.

Σ is the true covariance matrix, mI the scalar multiple of the identity closest to Σ , and $\tilde{\Sigma}$ the sample covariance matrix. r_1 , r_2 and d denote the distances between these three matrices (see Theorem 5). The errors on mI and $\tilde{\Sigma}$ are orthogonal by Theorem 3. $\hat{\Sigma}$ is the weighted average of mI and $\tilde{\Sigma}$ with minimum mean squared error. It is the orthogonal projection of Σ onto the line between mI and $\tilde{\Sigma}$.

Bayesian Interpretation

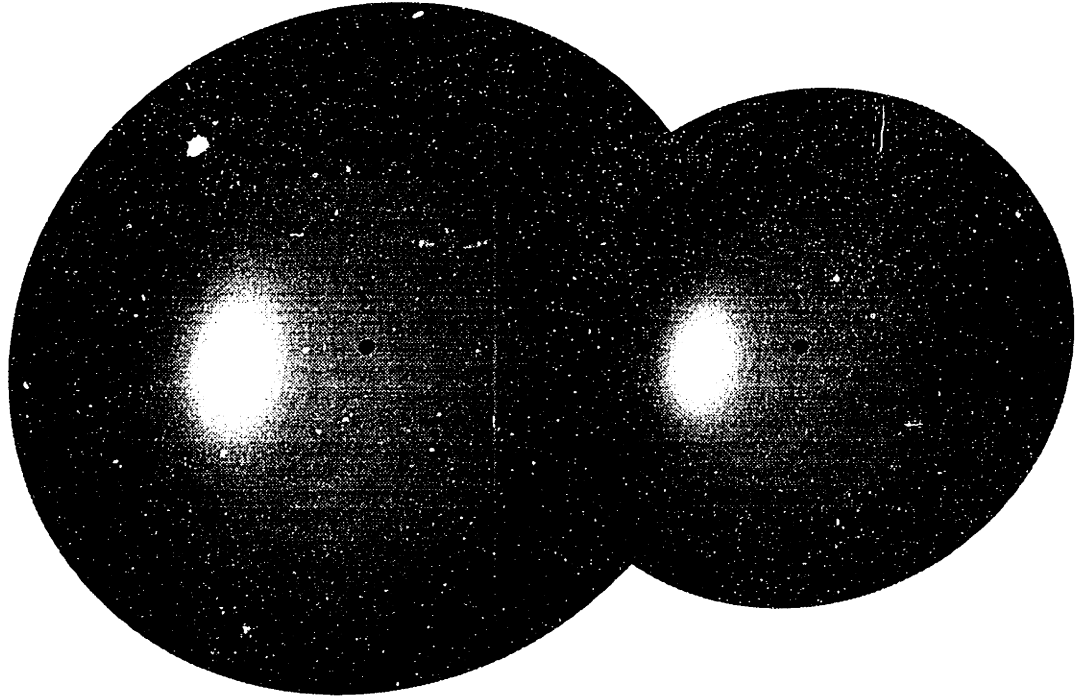


Figure F-3: Bayesian Interpretation.

The left sphere has center $\bar{\Sigma}$ and radius \hat{r}_1 . The right sphere has center $\tilde{\Sigma}$ and radius \hat{r}_2 . The distance between sphere centers is \hat{D} . If all we knew was that the true covariance matrix Σ lies on the left sphere, our best guess would be its center: the structured estimator $\bar{\Sigma}$. If all we knew was that the true covariance matrix Σ lies on the right sphere, our best guess would be its center: the sample covariance matrix $\tilde{\Sigma}$. Putting together both pieces of information, the true covariance matrix Σ must lie on the circle where the two spheres intersect, therefore our best guess is its center: the improved estimator $\hat{\hat{\Sigma}}$.

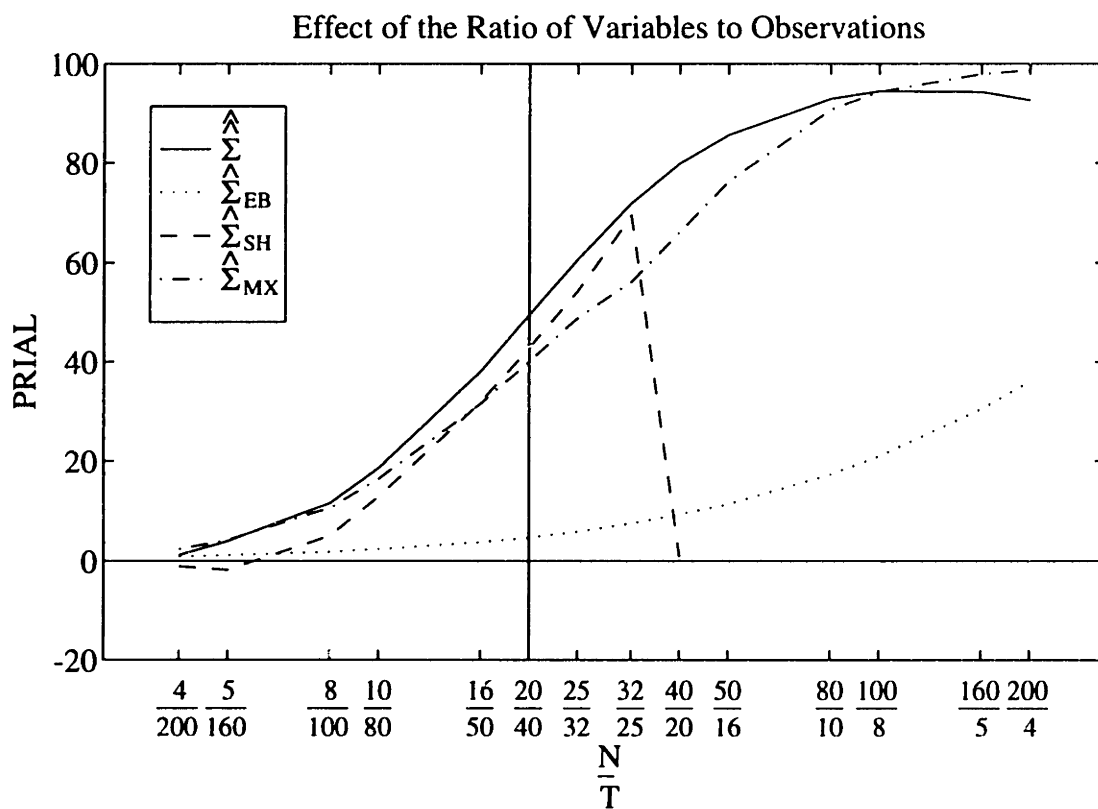


Figure F-4: Effect of the Ratio of Number of Variables to Number of Observations on the Percentage Relative Improvement in Average Loss (PRIAL). Estimators and parameters are described in Section 1.4.1. Based on 1,000 Monte-Carlo simulations. $\hat{\Sigma}_{SH}$ is not defined when $N/T > 2$ because the isotonic regression does not converge.

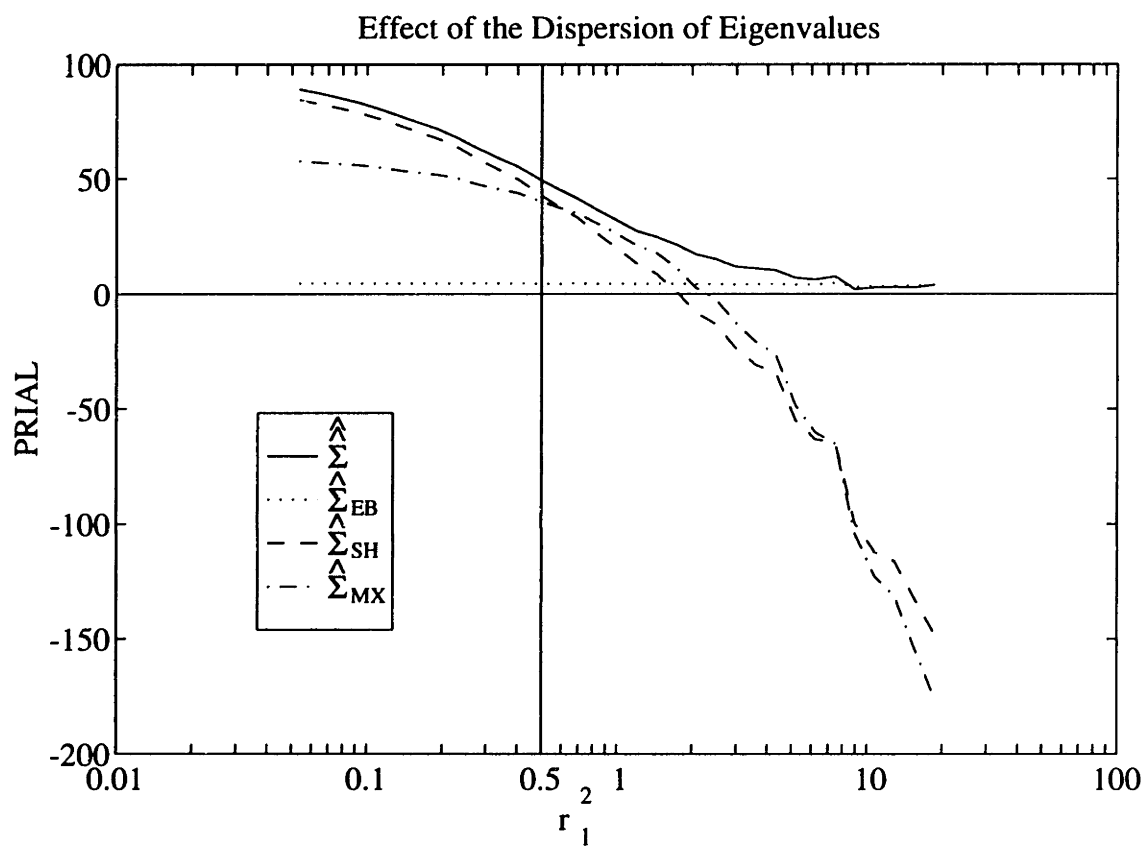


Figure F-5: Effect of the Dispersion of Eigenvalues on the Percentage Relative Improvement in Average Loss (PRIAL).

Estimators and parameters are described in Section 1.4.1. Based on 1,000 Monte-Carlo simulations.

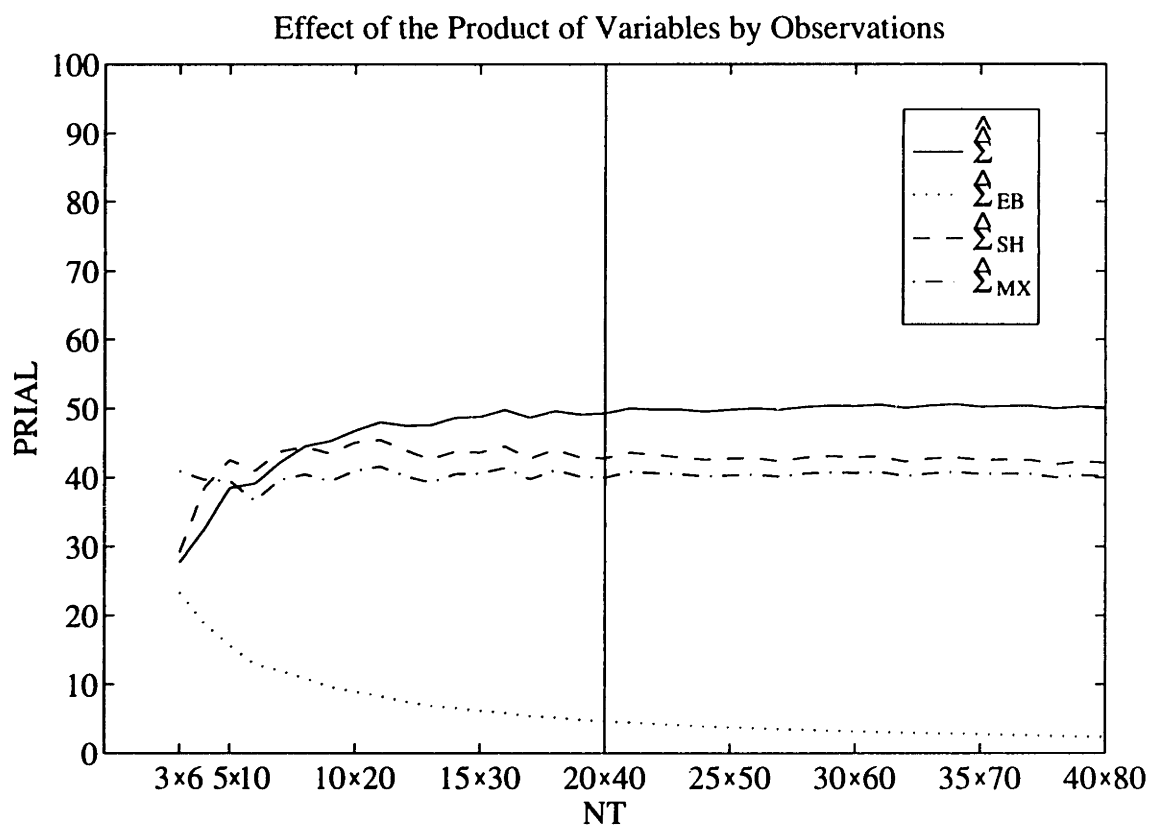


Figure F-6: Effect of the Product of Variables by Observations on the Percentage Relative Improvement in Average Loss (PRIAL).
 Estimators and parameters are described in Section 1.4.1. Based on 1,000 Monte-Carlo simulations.

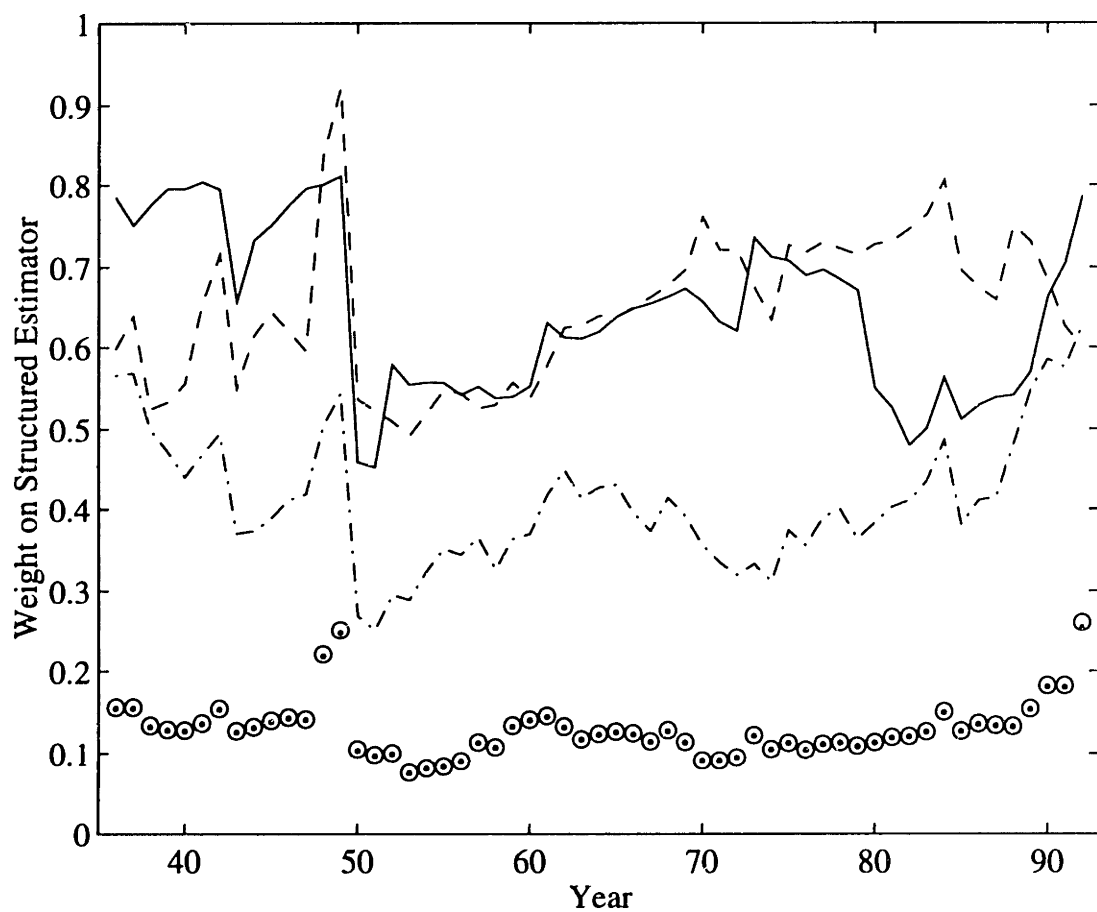


Figure F-7: Weights on Structured Estimators.

These weights are equal to $(\hat{r}_2^2 - \hat{\varphi})/\hat{D}^2$, see Theorem 10. Dots correspond to the structured estimator $\bar{\Sigma} = \hat{m}I$; circles, to the structured estimator of Appendix B.2; the dashed-dotted line, to Appendix B.1; the dashed line, to Appendix B.3; and the solid line, to Appendix B.4.

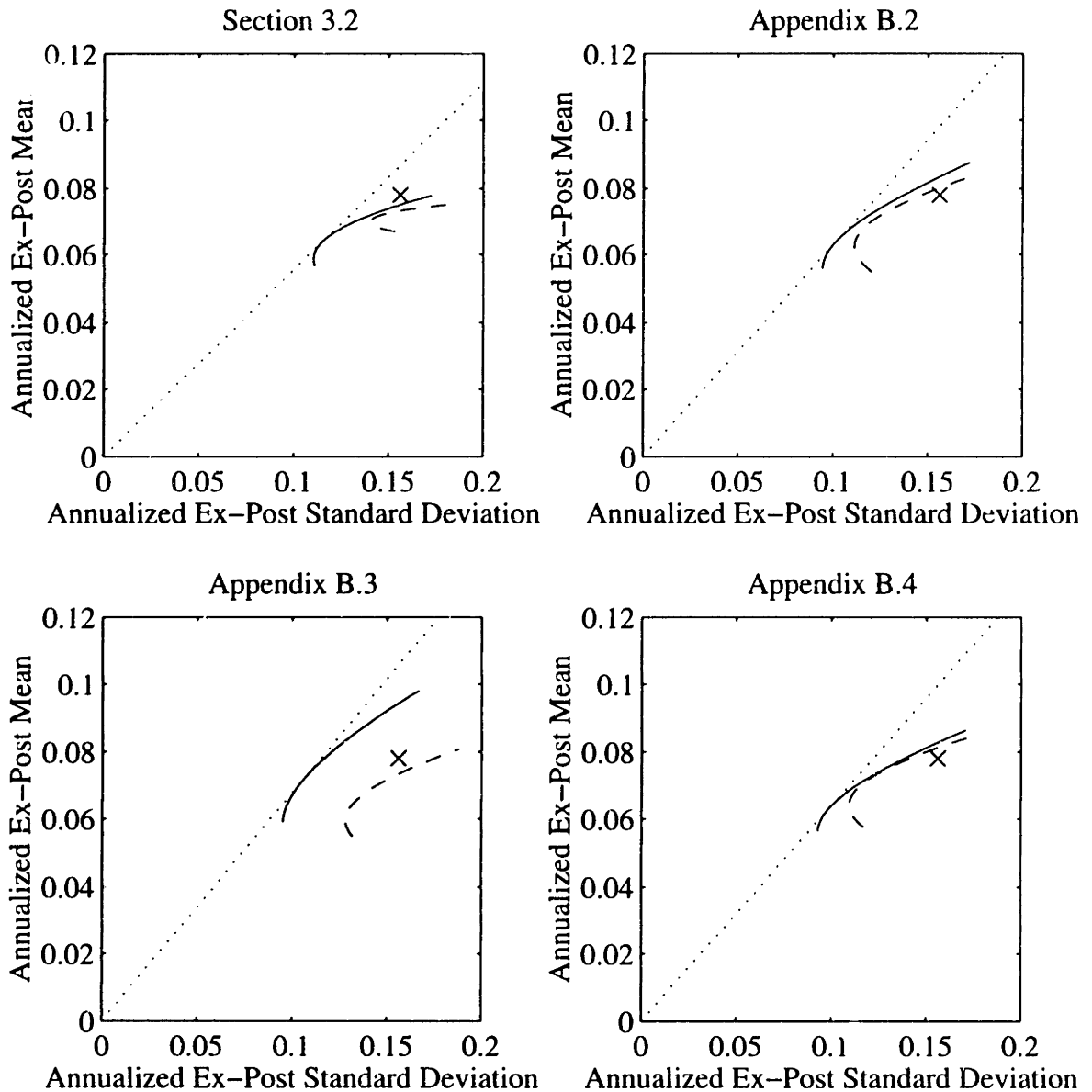


Figure F-8: Ex-Post Characteristics of Ex-Ante Constrained Minimum Variance Portfolios. Portfolios are constrained to have a specified beta between zero and one, and size zero. On each graph, portfolios obtained from a structured estimator are plotted as a dashed line, together with portfolios from the corresponding shrinkage estimator as a solid line. The title of each graph gives the section where the structured estimator is described. In the interest of space, the graph corresponding to Appendix B.1 is not shown. It closely resembles the one corresponding to Section 1.3.2. The symbol \times represents the CRSP value-weighted index, for reference. Shrinkage improves the risk-return tradeoff, moderately for the graphs on the left, and very slightly for the ones on the right.

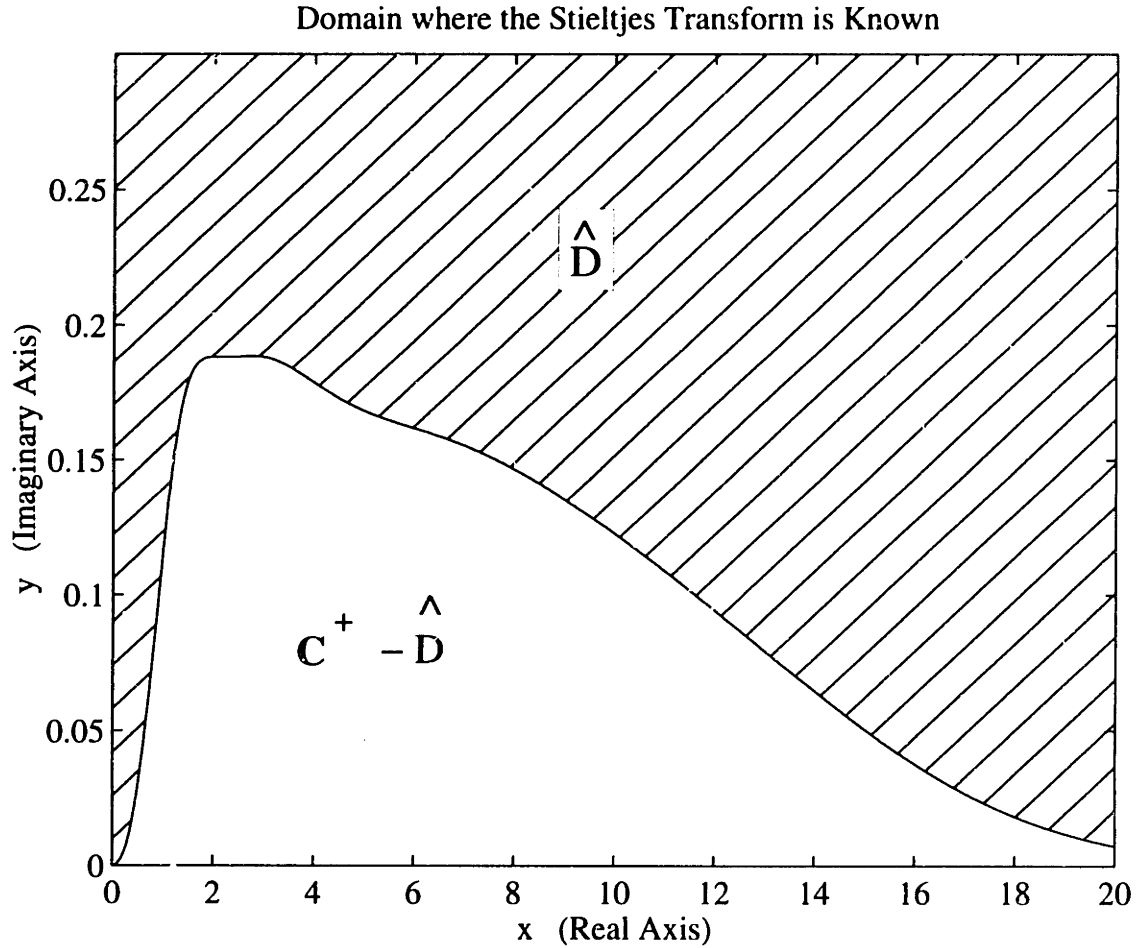


Figure F-9: Domain where the Value of $s_{L\hat{H}}$ is Known from Equation (A.4). The hatched domain represents a typical domain \hat{D} , cf. Appendix A. \hat{D} is the domain where an estimate $s_{L\hat{H}}$ of the Stieltjes transform of the true spectral c.d.f. H is known from Equation (A.4). The value of $s_{L\hat{H}}$ is not shown in this figure. The Stieltjes inversion formula ties the density $h(x)$ of true eigenvalues to the imaginary part of $s_{L\hat{H}}(x + i\varepsilon)$ for small $\varepsilon > 0$. Therefore we must extend $\text{Im}[s_{L\hat{H}}]$ from the hatched domain \hat{D} towards the real line. It means solving a Laplace equation with free boundary. This is an ill-posed problem. The degree of ill-posedness is proportional to how far the hatched domain is from the real line. In this simulation, ill-posedness is less severe around large eigenvalues (large x) than small ones (small x). This figure is generated from $T = 1000$ observations on $N = 100$ variables. The true spectral c.d.f. is the standard lognormal distribution. It has many small, clustered eigenvalues and a few large, more isolated ones. This is the same general shape as the eigenvalues of the covariance matrix of the returns on all stocks traded in the stock market.

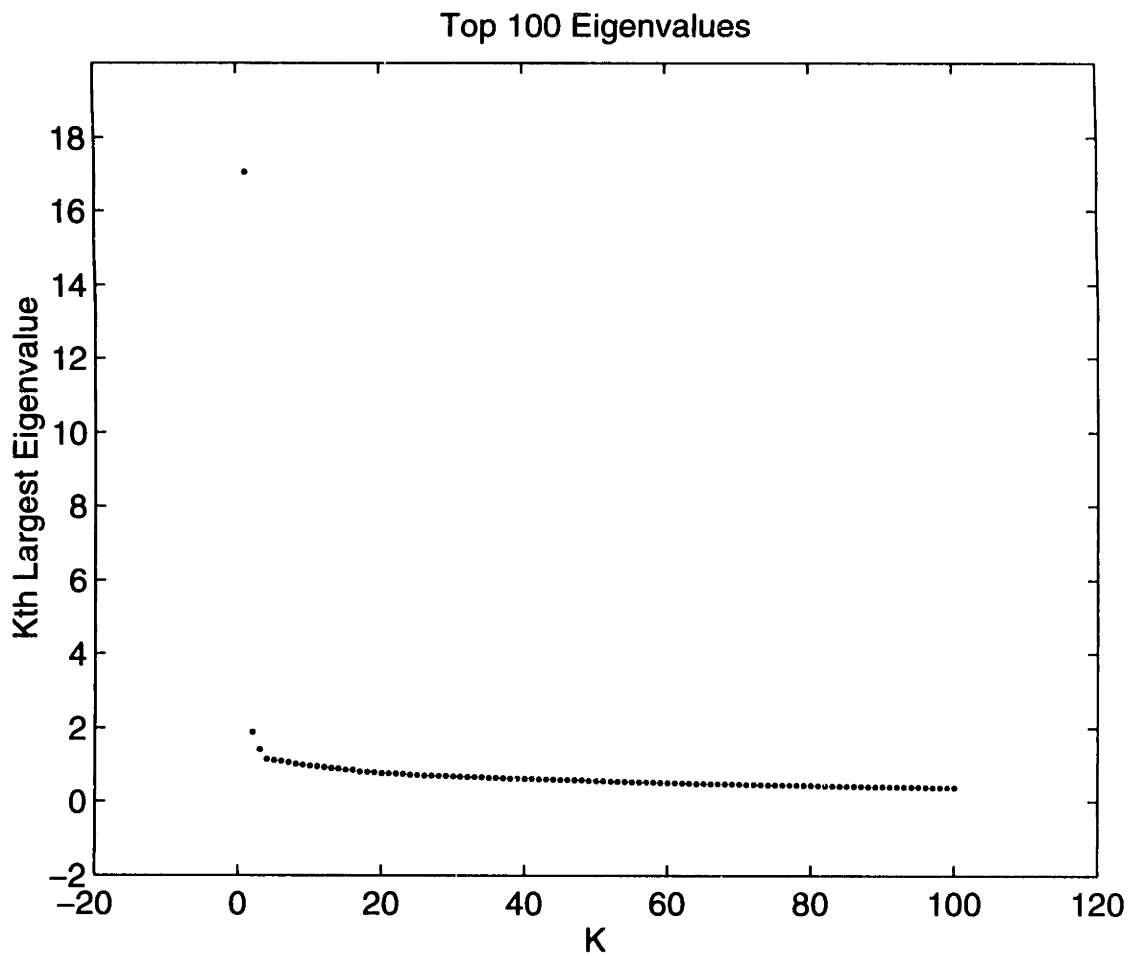


Figure F-10: Top 100 Eigenvalues of the Covariance Matrix of Stock Returns. The covariance matrix of NYSE and AMEX stock returns is estimated from daily CRSP data over 7/62-6/82. There are 5017 observations on 1019 stocks.

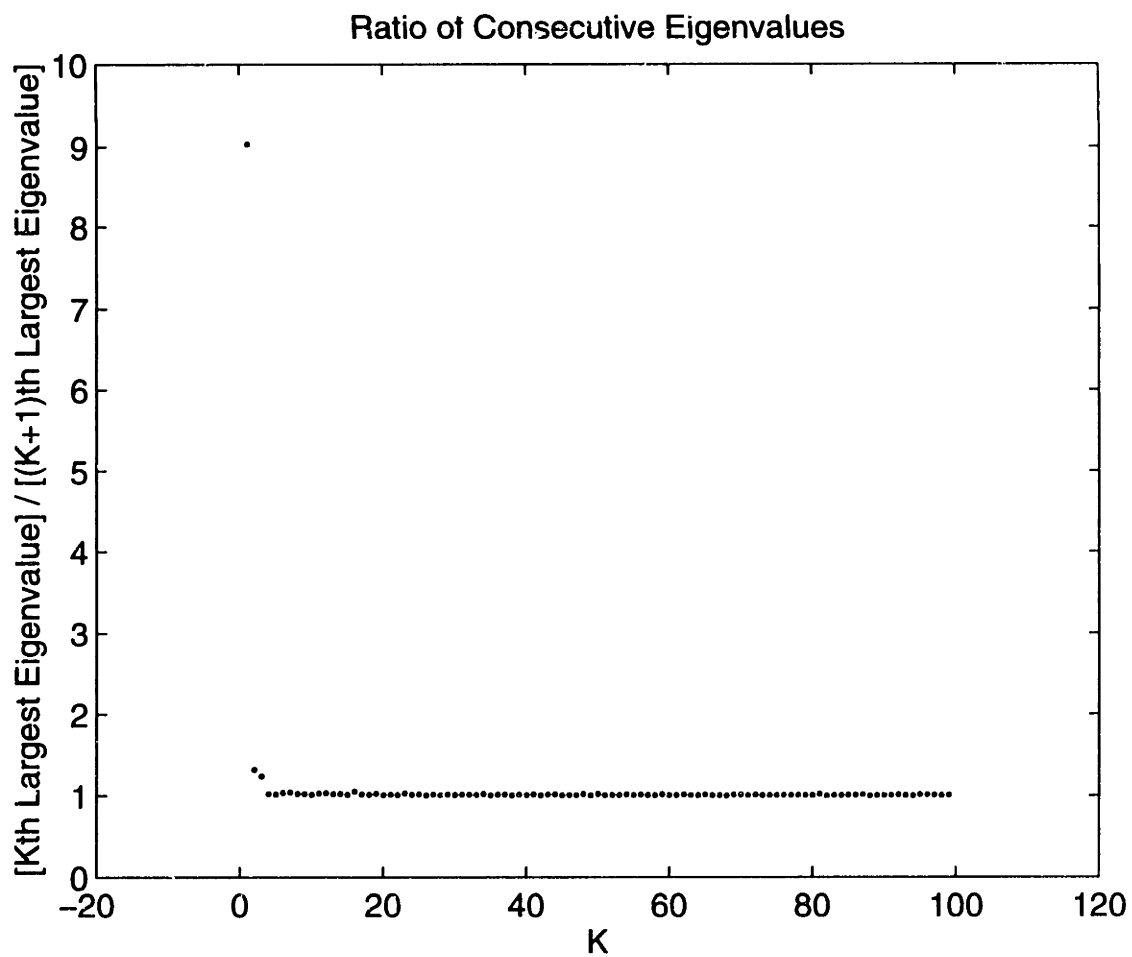


Figure F-11: Ratio of Consecutive Eigenvalues of the Covariance Matrix of Stock Returns. The ratio plots well above one for the first three eigenvalues only. These are the only eigenvalues between which a gap is apparent.

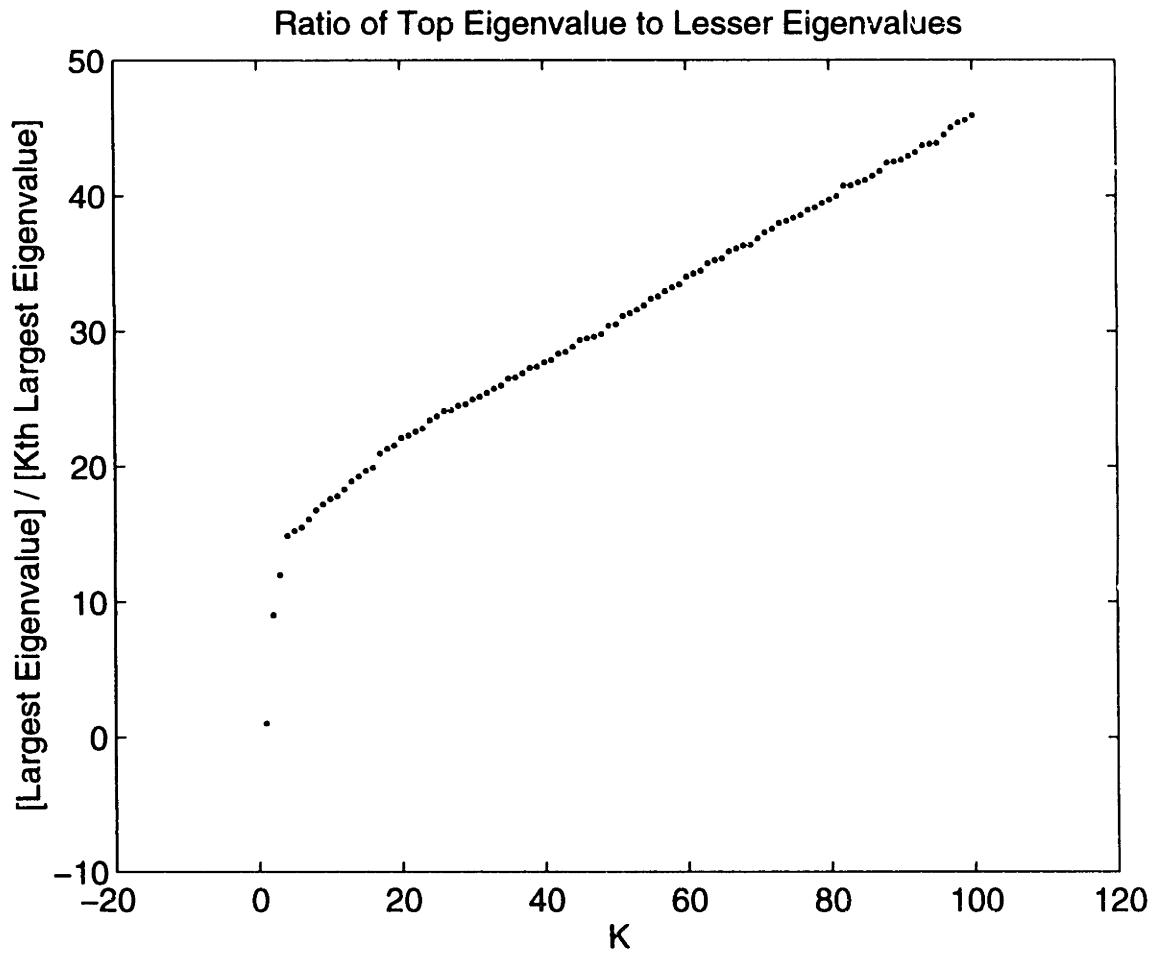


Figure F-12: Ratio of the Top Eigenvalue to Lesser Eigenvalues.
 The further we look down the ranking of eigenvalues, the more negligible they become with respect to the first one. In order to obtain negligible residuals, we must take a large number K of factors. This contradicts the requirement of a large gap between factors and residuals.

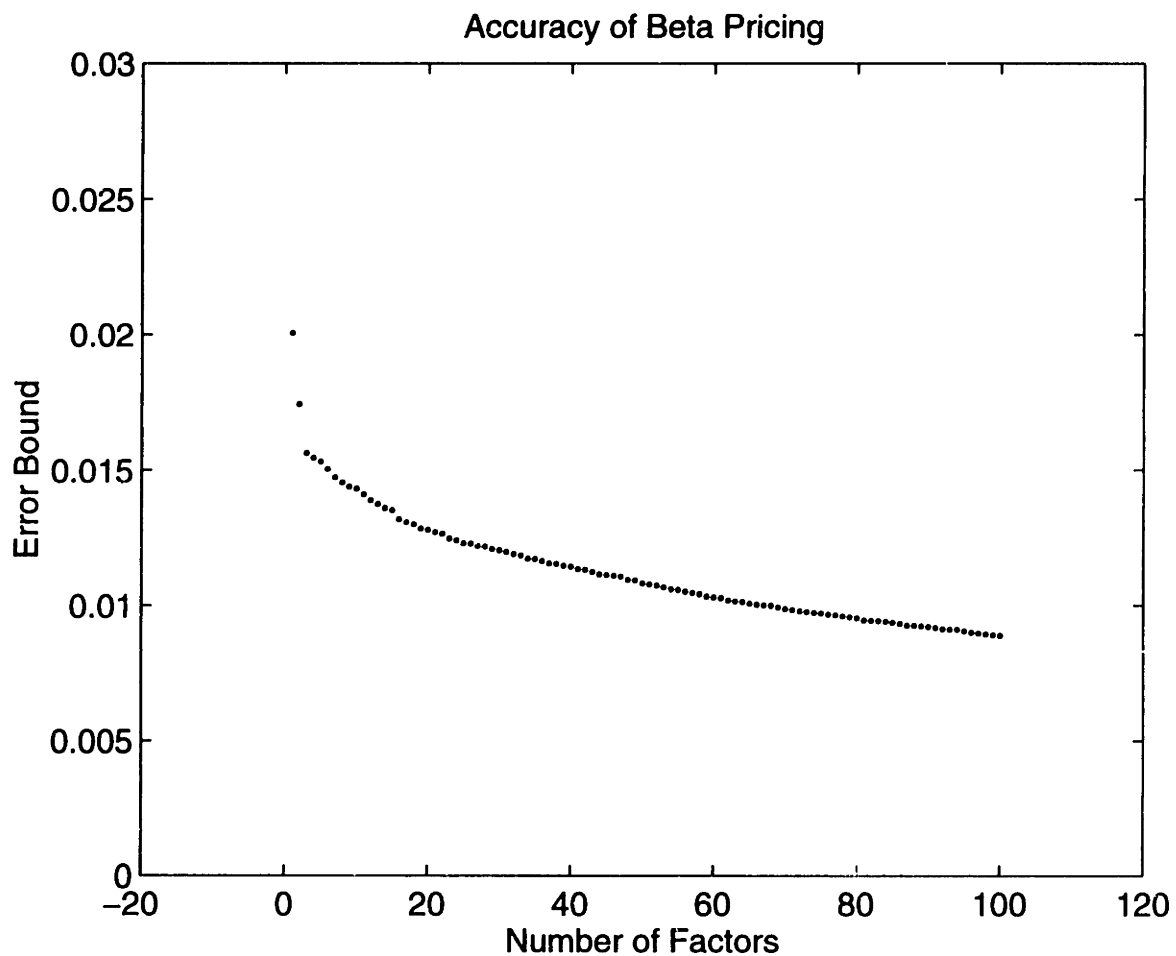


Figure F-13: Beta Pricing Error Bound.
This graph plots the upper bound on mean squared deviations from beta pricing in Equation (2.4).

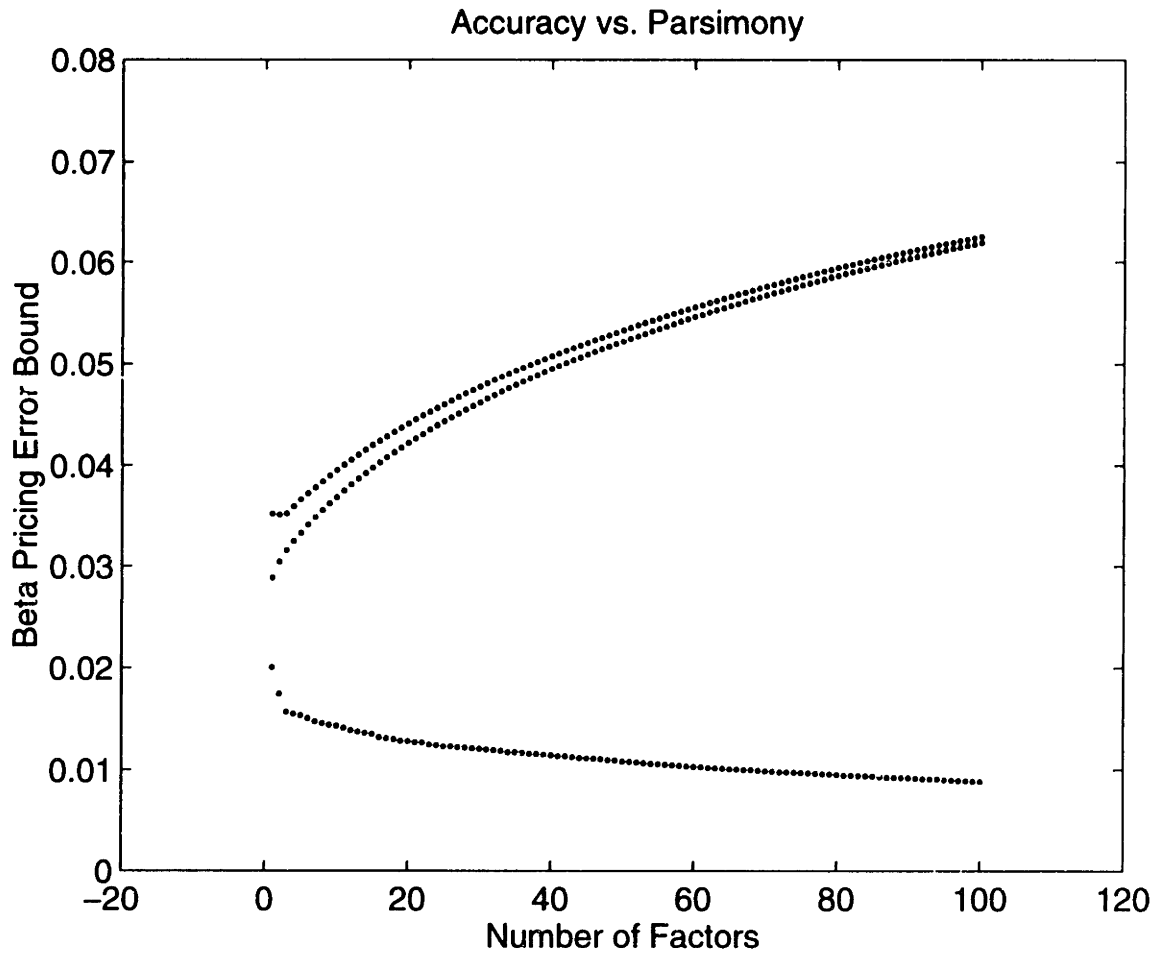


Figure F-14: Accuracy vs. Parsimony.

The lower dots plot the square root of the first term on the right hand side of Equation (2.7). It represents the deviation from beta pricing due to residual risk. It decreases in the number of factors K . The middle dots plot the square root of the second term on the right hand side of Equation (2.7). It represents the deviation from beta pricing due to risk premium estimation error. It increases in the number of factors K . The upper dots plot the square root of the sum of the two terms on the right hand side of Equation (2.7). It represents the total deviation from beta pricing. It is minimized for $K = 2$ factors. Choosing K to minimize the total deviation from beta pricing involves a trade-off between accuracy (with residual risk) and parsimony (with risk premium estimation error). The solution of this trade-off is the optimal number of factors in the δ -APT. The optimal number of factors is quite small. Even at the optimum, deviations from beta pricing are rather large.

Bibliography

- [1] Petr Adamek. Approximate factor structure: a test for number of factors. Technical report, MIT Sloan School of Management, 1994.
- [2] Yakov Amihud, Bent Jesper Christensen, and Haim Mendelson. Further evidence on the risk-return relationship. Technical report, New York University, 1992.
- [3] V. S. Bawa, Stephen J. Brown, and R. W. Klein. *Estimation Risk and Optimal Portfolio Choice*. Bell Laboratory Series. North Holland, New York, 1979. Studies in Bayesian Econometrics.
- [4] Fischer Black and Myron Scholes. The pricing of options and corporate liabilities. *Journal of Polititcal Economy*, 81:637–654, May-June 1973.
- [5] Stephen J. Brown. The number of factors in security returns. *Journal of Finance*, 44(5):1247–1262, December 1989.
- [6] Gary Chamberlain and Michael Rothschild. Arbitrage, factor structure, and mean-variance analysis on large asset markets. *Econometrica*, 51(5):1281–1304, September 1983.
- [7] Nai-fu Chen, Richard Roll, and Stephen A. Ross. Economic forces and the stock market. *Journal of Business*, 59:383–403, July 1986.
- [8] Gregory Connor and Robert A. Korajczyk. The arbitrage pricing theory and multifactor models of asset returns. In *Finance Handbook*, 1992. R. Jarrow, V. Maksimovic and W. Ziemba, eds.

- [9] Kent Daniel and Sheridan Titman. Evidence on the characteristics of cross-sectional variation in stock returns. Graduate School of Business, University of Chicago, March 1995.
- [10] D. K. Dey and C. Srinivasan. Estimation of a covariance matrix under Stein's loss. *Annals of Statistics*, 13(4):1581–1591, 1985.
- [11] Eugene F. Fama. *Foundations of Finance*. Basic Books, New York, 1970.
- [12] Eugene F. Fama and Kenneth R. French. The cross-section of expected stock returns. *Journal of Finance*, 1992.
- [13] Peter A. Frost and James E. Savarino. An empirical Bayes approach to portfolio selection. *Journal of Financial and Quantitative Analysis*, 21(3):293–305, September 1986.
- [14] L. R. Haff. Empirical Bayes estimation of the multivariate normal covariance matrix. *Annals of Statistics*, 8:586–597, 1980.
- [15] L. R. Haff. Solutions of the Euler-Lagrange equations for certain multivariate normal estimation problems. Unpublished manuscript, 1982.
- [16] Gur Huberman. A simple approach to arbitrage pricing. *Journal of Economic Theory*, 28:183–191, 1982.
- [17] J. D. Jobson and Bob Korkie. Estimation for Markowitz efficient portfolios. *Journal of the American Statistical Association*, 75(371):544–554, September 1980. Applications Section.
- [18] Shmuel Kandel and Robert F. Stambaugh. Portfolio inefficiency and the cross-section of expected returns. Technical report, Wharton School, 1994.
- [19] Bob Korkie. Corrections for trading frictions in multivariate returns. *Journal of Finance*, 44(5):1421–1434, December 1989.

- [20] Josef Lakonishok and Alan C. Shapiro. Systematic risk, total risk and size as determinants of stock market returns. *Journal of Banking and Finance*, 10:115–132, 1986.
- [21] A. Craig MacKinlay. Distinguishing among asset pricing theories: An *ex ante* analysis. Wharton School, University of Pennsylvania, April 1993.
- [22] Harry Markowitz. Portfolio selection. *Journal of Finance*, 7(1):77–91, March 1952.
- [23] V. A. Marčenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the U.S.S.R. - Sbornik*, 1(4):457–483, 1967.
- [24] Richard O. Michaud. The Markowitz optimization enigma: is ‘optimized’ optimal? *Financial Analysts Journal*, pages 31–42, January-February 1989.
- [25] Robb J. Muirhead. Developments in eigenvalue estimation. *Advances in Multivariate Statistical Analysis*, pages 277–288, 1987.
- [26] Richard Roll. A critique of the asset pricing theory’s test; part I: on past and potential testability of the theory. *Journal of Financial Economics*, 4:129–176, 1977.
- [27] Stephen A. Ross. The arbitrage theory of capital asset pricing. *Journal of Economic Theory*, 13:341–360, 1976.
- [28] Stephen A. Ross. Mutual fund separation in financial theory — The separating distributions. *Journal of Economic Theory*, 17:254–286, 1978.
- [29] Jay Shanken. The arbitrage pricing theory: Is it testable? *Journal of Finance*, 1982.
- [30] Jay Shanken. Multi-beta CAPM or equilibrium-APT?: A reply. *Journal of Finance*, 40(4):1185–1196, September 1985.
- [31] Jay Shanken. Nonsynchronous data and the covariance-factor structure of returns. *Journal of Finance*, 42(2):221–231, June 1987.
- [32] Jay Shanken. The current state of the arbitrage pricing theory. *Journal of Finance*, 47(4):1569–1574, September 1992.

- [33] William F. Sharpe. A simplified model for portfolio analysis. *Management Science*, January 1963.
- [34] Yo Sheena and Akimichi Takemura. Inadmissibility of non-order-preserving orthogonally invariant estimators of the covariance matrix in the case of Stein's loss. *Journal of Multivariate Analysis*, 41:117–131, 1992.
- [35] Jack W. Silverstein. Strong convergence of the empirical distribution of eigenvalues of large dimensional random matrices. *Journal of Multivariate Analysis*, 1994. Submitted.
- [36] Jack W. Silverstein and Sang-Il Choi. Analysis of the limiting spectral distribution of large dimensional random matrices. *SIAM Journal on Mathematical Analysis*, 1994. Submitted.
- [37] Jack W. Silverstein and Patrick L. Combettes. Signal detection via spectral theory of large dimensional random matrices. *IEEE Transactions on Signal Processing*, 40:2100–2105, 1992.
- [38] Charles Stein. Estimation of a covariance matrix. Rietz Lecture, 39th Annual Meeting IMS. Atlanta, GA., 1975.
- [39] Charles Stein. Series of lectures given at the University of Washington, Seattle, 1982.
- [40] Seha M. Tinic and Richard R. West. Risk and return: January vs. the rest of the yer. *Journal of Financial Economics*, 13:561–574, December 1984.
- [41] Kenneth W. Wachter. Probability plotting points for principal components. In *Proceedings of the Ninth Interface Symposium on Computer Science and Statistics*, pages 299–308, 1976. Hoaglin and Welsch, eds.
- [42] Y. Q. Yin. Limiting spectral distribution for a class of random matrices. *Journal of Multivariate Analysis*, 20:50–68, 1986.

Biographical Note

The author received a Bachelor's Degree in General Engineering (Diplôme d'Ingénieur) from the Ecole Polytechnique, Paris, France, in 1990. He obtained a Master's Degree in Statistics and Economics (Diplôme de Statisticien-Economiste) from the Ecole Nationale de la Statistique et de l'Administration Economique (ENSAE-SEA), Paris, France, in 1992.