

Multivariate Cointegration in Statistical Arbitrage

João Ramos Jungblut^{a*}

joaoramosjungblut@hotmail.com

^a*Faculty of Economics Sciences, Federal University of Rio Grande do Sul, João Pessoa, 52, Porto Alegre, Brazil, zip 90040-000*

Abstract

The present study investigates financial arbitrage strategies, overcoming limitations of techniques such as pairs trading. Through multivariate cointegration in the Johansen sense, the aim is to construct a long-short portfolio by identifying an equilibrium relationship among multiple assets. Parameters were estimated using maximum likelihood, while arbitrage opportunities were discovered by modeling the spread, targeting mean-reversion trades. This approach is consistent with economic theory and allows for determining portfolio asset weights. The empirical analysis collected data from Economatica, covering 203 stocks from the Bovespa index from 1997 to 2023. Results were compared to a benchmark and subjected to robustness tests, confirming the strategy's effectiveness in generating returns without significantly increasing risk.

Keywords: Statistical arbitrage; long-short portfolio; multivariate cointegration; time-series analysis; econometrics.

1 Introduction

Efficient markets are those in which asset prices fully reflect all available information [Fama, 1970]. This implies that profit opportunities are short-lived because participants quickly exploit any discrepancies, thereby restoring price equilibrium. Arbitrage, on the other hand, refers to the act of taking advantage of discrepancies in asset prices across different scenarios to gain risk-free profit. The goal is to buy an asset at a lower price in one market and sell it at a higher price in another, capitalizing on this difference. In essence, arbitrage relies on the concept of market inefficiency. The challenge is to identify inefficiencies that can be exploited to generate profit without incurring greater risks. In practice, this pursuit becomes feasible when analyzing

*The author expresses gratitude for the financial support from the Coordination for the Improvement of Higher Education Personnel (CAPES) through a scholarship.

various assets, precisely due to the complexity of modeling hidden patterns across multiple variables.

Pairs Trading is a statistical approach developed to identify pairs of assets with price co-movement, aiming to anticipate convergence to equilibrium in situations of overvaluation or undervaluation. This strategy was developed in the 1980s, according to [Cavalcante et al. \[2016\]](#), and gained prominence in academia after the publication of the distance method by [Gatev et al. \[2006\]](#), which proposes selecting pairs by minimizing a distance criterion (Euclidean). [Perlin \[2009\]](#) assesses the effectiveness of this approach in the Brazilian market, which I will explore in this work. However, [Taylor \[2011\]](#) argues that the distribution of stocks does not follow a normal distribution. Since the distance method is based on metrics such as Euclidean distance or correlation, which rely on the assumption of a normal distribution to capture short-term linear relationships, its reliability is compromised by the inherent inconsistency of such relations. Hence, it is more appropriate to use methods that capture non-linear dependence. For example, the use of copulas for constructing a mispricing index, as explored by [Xie et al. \[2016\]](#), and more recently, the mixture of copulas proposed by [Sabino da Silva et al. \[2023\]](#).

Cointegration, despite being a linear statistical method, is more appealing to work with because it captures long-term equilibrium relationships. This enables the identification of price divergences that align with a theory consistently. Some works provide practical insights, such as, for example, [Chan \[2021\]](#) and [Diamond \[2014\]](#). In this approach, we are essentially interested in modeling two assets to produce a stationary time series called spread. Thus, if the two assets have an equilibrium relation, this series tends to revert to the mean in the long term, allowing for long and short operations. An empirical analysis of the Brazilian market involving cointegration has already been conducted by [Caldeira and Moura \[2013\]](#). They found high profitability in applying cointegration but did not impose any restrictions on the pair selection universe.

Seeking to enhance the method, [Ramos-Requena et al. \[2017\]](#) introduced the use of the Hurst exponent to assess the quality of random walk diffusion, observing whether the spread is truly mean-reverting. This could be done through the variance test proposed by [Lo and MacKinlay \[1989\]](#), aiming to identify the hypothesis that financial asset prices follow or not a random walk. Alternatively, the evaluation could be performed through the half-life, calculated from the Ornstein-Uhlenbeck process, as employed by [Teixeira \[2014\]](#). The latter used it as a criterion to determine the exit time of operations to reduce losses and increase profitability. Finally, [Sarmiento and Horta \[2021\]](#) combines the techniques described above with the use of unsupervised machine learning models to improve efficiency in identifying pairs of ETFs in intraday data. The author employs the OPTICS and DBSCAN algorithms to narrow the

search for cointegrated pairs within specific clusters. Subsequently, filters are developed using the Hurst exponent and half-life for the selection of the most promising pairs.

That being said, some considerations must be made. The first of these concerns the sensitivity of these models to specific market conditions. The performance of pair-based strategies may be affected by extraordinary events or periods of extreme volatility, which can compromise the robustness and reliability of the results. Regarding this, [Palomar \[2020\]](#) demonstrates how to use the Kalman filter to make corrections to the spread when structural breaks occur, as well as how to weigh each asset in the operations.

Another significant shortcoming of the pair selection technique is the fact that it does not provide ways to construct a well-optimized portfolio. Simply identifying cointegrated pairs does not take into account effective diversification and resource allocation. In this regard, we built a cointegrated portfolio using [Johansen \[1988\]](#) test which allows for exploring cointegration in a multivariate case. Unlike traditional methods focusing solely on pairs, this methodology conducts cointegration tests across multiple assets concurrently. By doing so, it facilitates portfolio construction rooted in the principles of mean-reversion trading, leveraging spread analysis to identify and exploit price divergences from equilibrium points.

We used stock price data comprising the Bovespa index from 1997 to 2023, collected from Economica. The database was consistently divided into one-year in-sample periods and six-month out-of-sample periods. The study was conducted on the Brazilian market, an emerging market with the potential for identifying arbitrage opportunities compared to more established markets. This market also lacks sufficient research, especially regarding pairs trading strategies, such as [Perlin \[2009\]](#) and [Caldeira and Moura \[2013\]](#), stated earlier, which did not construct a well-optimized portfolio. The chosen study period coincides with the implementation of the Plano Real in Brazil. Also, this period encompasses various economic crises, providing a rich dataset to analyze how the strategy performed in the face of market volatility and economic uncertainty.

Backtesting unveils promising outcomes, showcasing the potential of statistical arbitrage to yield substantial rewards with minimal risk. Impressively, during the study period, this strategy not only surpassed benchmarks but also demonstrated lower drawdowns. By examining its performance over the selected period and addressing the research gap, valuable insights into the strategy's robustness and effectiveness across different market conditions can be obtained, offering implications for investors.

The structure of this work is organized as follows: in the Background section [2](#), we provide detailed explanations of the Vector Autoregressive (VAR), cointegration, Vector Error Correc-

tion Model (VECM), Johansen’s approach, and Maximum Likelihood Estimation (MLE). In the Methodology section 3, we present the Trading Strategy and Performance Methods. In the Empirical Analysis 4, we present the data, discuss the results, and validate them through the robustness test. Finally, in the Conclusion 5, we review the obtained results and fundamental concepts of the model in order to discuss potential improvements and future research directions.

2 Background

The models for implementing the trading strategy will be presented here. First, we will explain the VAR, the basis for a deeper understanding of the others. Following that, we will address the concept of cointegration and formulate the VECM. The Johansen test will be used to identify cointegrated vectors, and parameter estimation will then be performed through MLE. The notations used are based on the works of Bueno [2018].

2.1 VAR

The VAR expresses entire economic models by providing constraints and equations, allowing parameters to be estimated. It is significant because it examines the trajectory of endogenous variables in the presence of a structural shock. Following this approach, it is possible to use the information presented to separate long-term patterns from short-term ones, thereby identifying mispriced assets. These variations are studied using residuals caused by noise in the price series of financial assets, and modeling them leads to the VECM.

In a general sense, we can express a model of order p , with n endogenous variables interconnected by a matrix A :

$$AX_t = B_0 + \sum_{i=1}^p B_i X_{t-i} + B\epsilon_t, \quad \epsilon_t \sim i.i.d.$$

A is a $n \times n$ matrix that defines the simultaneous constraints among the variables forming the $n \times 1$ endogenous variables at time t , X_t . The constant vector B_0 is $n \times 1$, B_i is a $n \times n$ matrix, B is a diagonal matrix of standard deviations, and ϵ_t is a $n \times 1$ vector of uncorrelated random disturbances.

Examining the explanation for a bivariate situation helps better understand the model’s endogeneity. Consider these equations:

$$\begin{aligned} x_{1,t} &= b_{10} \boxed{-a_{12}x_{2,t}} + b_{11}x_{1,t-1} + b_{12}x_{2,t-1} + \sigma_{x_1}\epsilon_{x_1,t} \\ x_{2,t} &= b_{20} \boxed{-a_{21}x_{1,t}} + b_{21}x_{1,t-1} + b_{22}x_{2,t-1} + \sigma_{x_2}\epsilon_{x_2,t}. \end{aligned}$$

In this context, (i) $x_{1,t}$ and $x_{2,t}$ are stationary, (ii) $\epsilon_{x_{1,t}} \sim RB(0,1)$ and $\epsilon_{x_{2,t}} \sim RB(0,1)$, and (iii) $\epsilon_{x_{1,t}} \perp \epsilon_{x_{2,t}} \Rightarrow Cov(\epsilon_{x_{1,t}}, \epsilon_{x_{2,t}}) = 0$.

These conditions are required for the model's validity and the accurate interpretation of the results. By ensuring the stationarity of variables and the independence of error terms, we may proceed with explaining the reduced form of the simple model, which becomes:

$$\begin{aligned} X_t &= A^{-1}B_0 + A^{-1} \sum_{i=1}^p B_i X_{t-i} + A^{-1}B\epsilon_t \\ &= \Phi_0 + \sum_{i=1}^p \Phi_i X_{t-i} + e_t. \end{aligned}$$

The simple model's reduced form is

$$X_t = \Phi_0 + \Phi_1 X_{t-1} + e_t,$$

where $\Phi_0 = A^{-1}B_0$, $\Phi_1 = A^{-1}B_1$, and $Ae_t = B\epsilon_t$. The stability condition is ensured when the eigenvalues are positioned outside the unit circle ($I - \Phi_1 L$), which ensures convergence over time. Moving on, the question arises about the presence of trends in the variables, i.e., if one may predict the other under what circumstances. This question introduces the notion of cointegration.

2.2 Cointegration

According to [Engle and Granger \[1987b\]](#), the elements of the vector X_t , $n \times 1$, are said to be cointegrated of order (d, b) , denoted by $X_t \sim CI(d, b)$, if (i) all elements of X_t are integrated of order d , $I(d)$, and (ii) there exists a non-zero vector β such that

$$u_t = X_t' \beta \sim I(d - b), \quad b > 0.$$

Therefore, when $X_t' \beta = 0$, β is the cointegration vector that defines a linear combination among the elements of X_t .

Consider $\beta = [\hat{\beta}_1, \hat{\beta}_2]$, which defines the long-term equilibrium between the variables $I(1)$, $x_{1,t}$, and $x_{2,t}$. Then,

$$\begin{bmatrix} x_{1,t} & x_{2,t} \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \hat{\beta}_1 x_{1,t} + \hat{\beta}_2 x_{2,t} = 0$$

$$\begin{bmatrix} x_{1,t} & x_{2,t} \end{bmatrix} \begin{bmatrix} 1 \\ \beta_2 \end{bmatrix} = x_{1,t} + \beta_2 x_{2,t} = 0.$$

With the normalization of $\beta_2 = \frac{\hat{\beta}_2}{\hat{\beta}_1}$, u_t can be considered the residual of one coordinate of X_t against the other variables. Thus, the variables are cointegrated, and β generates the residual whose order of integration is lower than the original variables. The same holds for cases with more than two variables.

The VAR is significant because it attempts to predict the trajectory of endogenous variables in the face of a structural shock. To trade, we must first assess the shock and determine entry and exit points. Given the condition for a stationary disturbance, the logical next step is to test the residuals, fit the best VAR model, and send the information to the VECM.

2.3 VECM

The VECM is nothing more than a conventional VAR, but it includes the error correction term. To visualize this, we need to consider a cointegration relationship given by:

$$x_{1,t} = \mu + \beta x_{2,t} + u_t.$$

So, there are ways to manipulate the VAR such that, if cointegration exists, the original model can be rewritten for the residuals to enter explicitly:

$$\begin{aligned} \Delta x_{1,t} &= \alpha_1 \hat{u}_{t-1} + \sum_{j=1}^{p-1} \lambda_{11,j+1} \Delta x_{1,t-j} + \sum_{j=1}^{p-1} \lambda_{12,j+1} \Delta x_{2,t-j} + e_{x_{1,t}} \\ \Delta x_{2,t} &= \alpha_2 \hat{u}_{t-1} + \sum_{j=1}^{p-1} \lambda_{21,j+1} \Delta x_{1,t-j} + \sum_{j=1}^{p-1} \lambda_{22,j+1} \Delta x_{2,t-j} + e_{x_{2,t}}. \end{aligned}$$

In the multivariate model, each X_t is a $n \times 1$ vector of endogenous variables.

Now, consider the VAR at the level, ignoring the presence of constants, to better comprehend the VECM.

$$\begin{aligned} X_t &= \Phi_1 X_{t-1} + \Phi_2 X_{t-2} + \dots + \Phi_p X_{t-p} + e_t \\ [I - (\Phi_1 L + \Phi_2 L^2 + \dots + \Phi_p L^p)] X_t &= e_t \\ \Phi(L) X_t &= e_t. \end{aligned}$$

Note, when $L = 1$,

$$\Phi(1) = [I - (\Phi_1 + \Phi_2 + \dots + \Phi_p)] = -\Phi.$$

The characteristic polynomial is given by:

$$\Phi(Z) = I - \sum_{i=1}^p \Phi_i Z^i,$$

where Z is a diagonal matrix of n elements. Linear algebra tells us that if the matrix's determinant is zero, its rank is not full. That is, $[\Phi(I)] = 0 \iff \text{rank}(\Phi) < n$. As a result, the process has a unit root, and Z can be factored as follows:

$$\Phi(Z) = (I - Z)(I - \lambda_1 Z)(I - \lambda_2 Z) \dots (I - \lambda_p Z).$$

Remember that matrix's rank is the number of independent rows or columns. The number of columns and rows will always be fewer than or equal to the rank. According to [Bueno \[2018\]](#), this allows us to state Granger's theorem, originally proposed by [Gonzalo and Granger \[1995\]](#).

Theorem 1 (Granger). *Let $\Phi : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}^{n \times n}$ be a linear map s.t.*

- (i) $|\Phi(Z)| = 0 \Rightarrow Z > I$, where I denotes the $n \times n$ identity matrix;
- (ii) $0 < \text{rank}(\Phi) = r < n$.

Then, $\exists \alpha \in \mathbb{R}^{n \times r}$ and $\beta \in \mathbb{R}^{r \times n}$ such that $\Phi = \alpha\beta$.

The theorem expresses the idea that Φ can be decomposed into two multiplicative matrices, β is the matrix of cointegration and α is the matrix of adjustment. From this, we derive the VECM, recursively adding and subtracting past terms to generalize the equation. Thus, let

$$X_t = \Phi X_{t-1} + \sum_{i=1}^2 \Lambda_i \Delta X_{t-i} + e_t,$$

with $\Lambda_i = \sum_{j=1+i}^3 \Phi_j$, $i = 1, 2$.

The model is called error correction because it explains ΔX_t by two components: the short-term factors and the long-term relationship given by the coordinates of the vector of endogenous variables. It becomes evident, therefore, that in the presence of cointegration, it is always possible to associate the VAR with error correction, and that is precisely what the representation theorem deals with, [Engle and Granger \[1987a\]](#).

Theorem 2 (Granger Representation Theorem). *If $X_t \sim CI(1, 1)$, X_t has a representation in the form of a VECM.*

2.4 Johansen's approach

Johansen presents a test to identify the rank of the matrix Φ and, consequently, estimate the cointegration vectors included in the matrix β . His methodology is intriguing because it is performed concurrently with the estimation of the cointegration model.

Consider the following non-stationary equation system in cointegration:

$$X_t = \beta_0 + \beta_1 X_{t-1} + u_t,$$

in vector form,

$$\begin{bmatrix} x_{1,t} \\ x_{2,t} \end{bmatrix} = \begin{bmatrix} \beta_{10} \\ \beta_{20} \end{bmatrix} + \begin{bmatrix} \beta_{11} & \alpha_{11} \\ \alpha_{21} & \beta_{21} \end{bmatrix} \begin{bmatrix} x_{1,t-1} \\ X_{2,t-1} \end{bmatrix} + \begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}.$$

In this system, x_{1t} and x_{2t} are non-stationary $I(1)$ processes. A linear combination exists that is $I(0)$ when they are cointegrated. We can express the VAR model in terms of the $I(0)$ variables only.

$$\begin{bmatrix} x_{1,t} - x_{1,t-1} \\ x_{2,t} - x_{2,t-1} \end{bmatrix} = \begin{bmatrix} \beta_{10} \\ \beta_{20} \end{bmatrix} + \begin{bmatrix} \beta_{11} - 1 & \alpha_{11} \\ \alpha_{21} & \beta_{21} - 1 \end{bmatrix} \begin{bmatrix} x_{1,t-1} \\ X_{2,t-1} \end{bmatrix} + \begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}.$$

$$\begin{bmatrix} \Delta x_{1,t} \\ \Delta x_{2,t} \end{bmatrix} = \begin{bmatrix} \beta_{10} \\ \beta_{20} \end{bmatrix} + \begin{bmatrix} \Phi_{11} & \Phi_{12} \\ \Phi_{21} & \Phi_{22} \end{bmatrix} \begin{bmatrix} x_{1,t-1} \\ X_{2,t-1} \end{bmatrix} + \begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}.$$

Now, the model is represented in the VECM form, as explained by the Granger Representation Theorem.

$$\Delta X_t = \beta_0 + \Phi X_{t-1} + u_t$$

The matrix's rank indicates the number of linearly independent rows and columns. If one of the rows cannot be stated as a multiple of the other, they are considered independent. As a result, the existence of a $I(0)$ linear combination for x_1 and x_2 depends on the rank. [Johansen \[1991\]](#) proposed two tests to evaluate the matrix's rank: the trace test and the maximum eigenvalue test. We are solely interested in the maximum eigenvalue test for this work, where the hypothesis is formulated as follows, given that r is the maximum number of cointegrated eigenvectors.

$$H_0 : \text{rank}(\Phi) < r \quad \text{or} \quad \Phi = \alpha\beta'$$

Even after defining the rank, the matrices α and β are not identifiable as they form an

overparameterization of the model. However, we can delimit the cointegration space to $span(\beta)$.

2.5 MLE

Parameters estimation can be done by maximum likelihood. Drawing from the discussions in [Maddala and Kim \[1998\]](#) and its references, consider the complete VAR:

$$X_t = \Phi_1 X_{t-1} + \Phi_2 X_{t-2} + \dots + \Phi_p X_{t-p} + \delta d_t + e_t.$$

This model can be represented in error correction form, and the objective is to estimate using Maximum Likelihood subject to the constraint of incomplete rank: $\Phi = \alpha\beta'$.

$$\Delta X_t = \alpha\beta' X_{t-1} + \sum_{i=1}^p \Lambda_i \Delta X_{t-i} + \delta d_t + e_t$$

Consider vectorizing the above model, where $\Upsilon_{0,t} = \Delta X_t$, $\Upsilon_{1,t} = X_{t-1}$, $\Upsilon_{2,t} = [\Delta X'_{t-1} \Delta X'_{t-2} \dots, d'_t]$, and $\Lambda = [\Lambda_1 \Lambda_2 \dots \delta]$. The VECM is simplified to:

$$\Upsilon_{0,t} = \alpha\beta' \Upsilon_{1,t} + \Lambda \Upsilon_{2,t} + e_t.$$

The problem is to maximize the likelihood function subject to nonlinear constraints on the parameters given by:

$$\ln L(\alpha, \beta, \Lambda, \Sigma) = -\frac{T}{2} \ln |\Sigma| - \frac{1}{2} \sum_{t=1}^T e'_t \Sigma^{-1} e_t,$$

with the first-order conditions being:

$$\begin{aligned} 0 &= \sum_{t=1}^T (\Upsilon_{0,t} - \alpha\beta' \Upsilon_{1,t} + \hat{\Lambda} \Upsilon_{2,t}) \Upsilon'_{2,t} \\ &\Rightarrow \sum_{t=1}^T \Upsilon_{0,t} \Upsilon'_{2,t} = \alpha\beta' \sum_{t=1}^T \Upsilon_{1,t} \Upsilon'_{2,t} + \hat{\Lambda} \sum_{t=1}^T \Upsilon_{2,t} \Upsilon'_{2,t}. \end{aligned}$$

Renaming the variables as $\Pi_{ij} = \frac{\sum_{t=1}^T \Upsilon_{ij} \Upsilon'_{ij}}{T}$, we can rewrite the previous equation as:

$$\Pi_{02} = \alpha\beta' \Pi_{12} + \hat{\Lambda} \Pi_{22}.$$

[Bueno \[2018\]](#) shows that to obtain e_t , we do not need $\hat{\Lambda}$. First, let's regress $\Upsilon_{0,t}$ on $\Upsilon_{2,t}$, that is,

$$\Upsilon_{0,t} = B \Upsilon_{2,t} + r_{0,t}$$

so that $\hat{B} = (\Upsilon_{0,t}\Upsilon'_{2,t})(\Upsilon_{2,t}\Upsilon'_{2,t})^{-1} = \Pi_{02}\Pi_{22}^{-1}$, and obtain the residuals

$$r_{\hat{0},t} = \Upsilon_{0,t} - \Pi_{02}\Pi_{22}^{-1}\Upsilon_{2,t}.$$

Next, regressing $\Upsilon_{1,t}$ on $\Upsilon_{2,t}$ and obtain the residuals $r_{\hat{1},t} = \Upsilon_{1,t} - \Pi_{12}\Pi_{22}^{-1}\Upsilon_{2,t}$. Note that:

$$\hat{e}_t = r_{\hat{0},t} - \alpha\beta r_{\hat{1},t}.$$

Thus, the quantities $r_{\hat{0},t}$ and $r_{\hat{1},t}$ can be obtained from the auxiliary regressions, and we can estimate α and β .

The new function to maximize becomes:

$$\ln L(\alpha, \beta, \Sigma) = -\frac{T}{2}\ln|\Sigma| - \frac{1}{2}\sum_{t=1}^T (r_{\hat{0},t} - \alpha\beta r_{\hat{1},t})'\Sigma^{-1}(r_{\hat{0},t} - \alpha\beta r_{\hat{1},t}).$$

One way is to estimate α and Σ for a given β ,

$$\hat{\alpha}(\beta) = S_{01}\beta(\beta'S_{11}\beta)$$

$$\hat{\Sigma}(\beta) = S_{00} - S_{01}\beta(\beta'S_{11}\beta)^{-1}\beta'S_{10}$$

with $S = T^{-1}\sum_{t=1}^T \hat{r}_{i,t}\hat{r}'_{j,t}$, $j = 0, 1$.

The matrix β has not been estimated yet; nonetheless, starting from the last derivation, β can be found by maximizing the likelihood.

$$L(\beta)^{-\frac{2}{T}} = |\hat{\Sigma}(\beta)|,$$

where according to [Johansen \[1991\]](#), we can obtain it through the maximum eigenvalue test.

$$L_{max}^{-\frac{2}{T}} = |S_{00}| \prod_{i=1}^n (1 - \hat{\lambda}_i).$$

Thus, the cointegrated eigenvectors will define the matrix

$$\hat{\beta} = [\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_r]$$

which is an $n \times r$ matrix.

In this context, let's clarify some details further.

$$\Pi_{02} = \alpha\beta'\Pi_{12} + \hat{\Lambda}\Pi_{22}$$

contains the coefficients associated with cointegration, with α and β being the coefficients we want to estimate. $r_{0,t}$ and $r_{1,t}$ are related to the variables deviating from their cointegration relationship (observed divergences), where $r_{0,t}$ is the residual of X_t (level) and $r_{1,t}$ is the residual of ΔX_t (trend). Furthermore, as we are only looking at the cointegrated space, the residual presented will be

$$e_t = \alpha\beta'X_t.$$

3 Methodology

The methodology is separated into two parts. The trading strategy is introduced in the first phase, with details on how buy and sell signals are identified based on the spread and when trades are executed. The second section introduces the performance metrics that will be used to evaluate the strategy's performance later.

3.1 Trading Strategy

The chosen strategy is similar to pairs trading using the cointegration method, according to books like [Chan \[2021\]](#) and studies like [Caldeira and Moura \[2013\]](#) and [Diamond \[2014\]](#). On the other hand, our approach distinguishes itself by attempting to build a cointegrated portfolio. The central proposal is to identify periods when the movement of a financial asset deviates from the others in the short term, anticipating a long-term convergence to equilibrium.

Detecting these price divergences requires observing the residuals of the linear relationships. Although the strategy is almost identical to pairs trading, we have n financial assets in the multivariate model. As explained in [Section 2](#), the residuals are given by:

$$e_t = \Phi X_t.$$

We shall use a separate nomenclature here to clarify the strategy's implementation. During the out-of-sample phase, the computed residuals will be referred to as spread and represent the vector of price divergences to be exploited. The Z-score, which is effectively the standardized spread, is calculated from these residuals and is represented as $Z_{\text{score}} = \frac{e_t - \mu_e}{\sigma_e}$.

Meanwhile, the normalized vector $\Phi = [1, \frac{\hat{\Phi}_2}{\hat{\Phi}_1}, \frac{\hat{\Phi}_3}{\hat{\Phi}_1}, \dots, \frac{\hat{\Phi}_F}{\hat{\Phi}_1}]$, obtained by Maximum Likelihood

Estimation (MLE), will constitute the hedge ratios. To determine the portfolio weights, w , we divide the hedge ratios by their norm

$$w = \frac{\Phi}{\|\Phi\|}.$$

This way, we ensure that only available capital is used; the remaining question is whether this is the best way to weigh the portfolio once assets with large value parameters concentrate all capital. The portfolio return is calculated as follows:

$$R_{port} = wR_t.$$

The strategy is built on the mean-reversion principle. Entry points are defined by the spread's distance from the mean, whereas exit points are determined by the spread's proximity to the mean. The Z-score is important because it allows us to determine entry and exit locations based on standard deviations.

$$\left\{ \begin{array}{l} \text{if } Z_{score} \leq -2, \quad \text{then } S_t = 1 \text{ (open long position)} \\ \text{if } Z_{score} \geq 2, \quad \text{then } S_t = -1 \text{ (open short position)} \\ \text{if } Z_{score} = 0, \quad \text{then } S_t = 0 \text{ (close position)} \\ \text{if } Z_{score} \leq -3, \quad \text{then } S_t = 2 \text{ (stop loss)} \\ \text{if } Z_{score} \geq 3, \quad \text{then } S_t = -2 \text{ (stop loss)} \end{array} \right.$$

Here, S_t denotes the trading signal at time t , where 1 indicates the initiation of a long position, -1 the initiation of a short position, 0 denotes the closing signal, and 2 or -2 implies the activation of a stop loss. Consequently, when the Z-score reaches -2 , a long position is opened in the portfolio ($R_{port} = wR_t$), closing it with a positive return at a Z-score of 0 or with a loss at -3 . Conversely, when the Z-score attains 2, a short position is opened ($R_{port} = -wR_t$), closing it at 0 or 3. A schematic representation in Figure 1 illustrates these events, wherein dashed green lines signify entry points, dashed black lines represent exit points, the dashed red line indicates the stop loss level, and the blue line depicts the spread.

3.2 Performance Measures

The evaluation of the investment portfolio's performance involves the computation of diverse measures. These metrics thoroughly evaluate the portfolio's risk and return attributes, facilitating a more nuanced understanding of its risk and return characteristics across different market conditions. The following key measures are calculated: annualized return, annualized standard deviation, annualized Sharpe ratio, annualized Sortino ratio, annualized Value

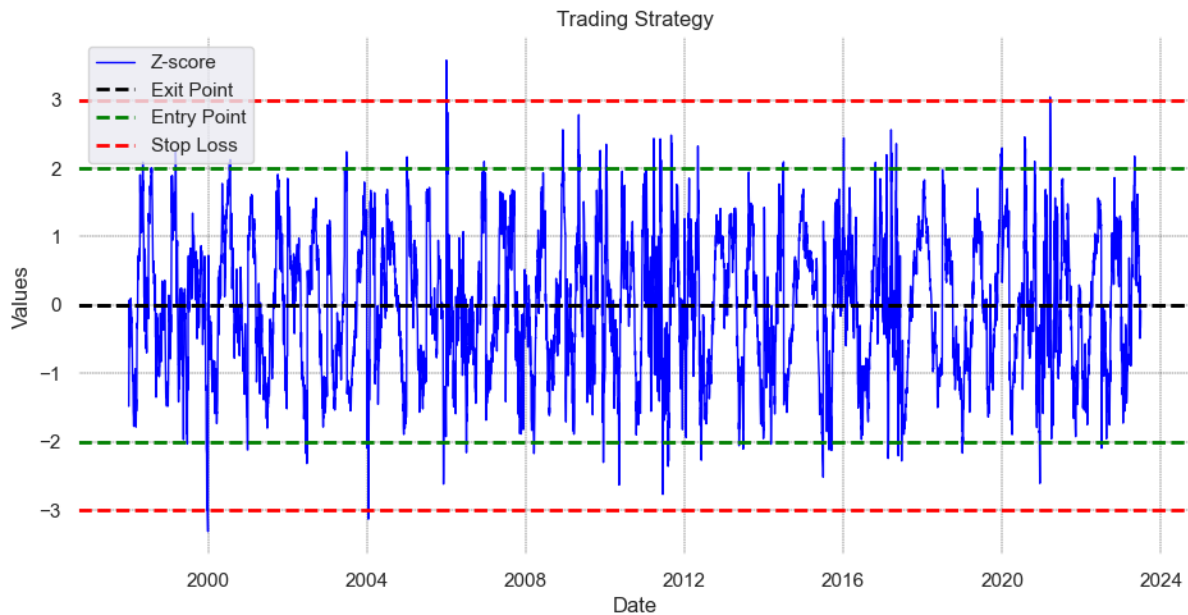


Figure 1: Graph depicting the spread of cointegrated portfolios from 1998 to 2023 for generating trading signals

at Risk (VaR), $VaR_{0.95}$, annualized Conditional Value at Risk (CVaR), $CVaR_{0.95}$, and Worst-Drawdown. These measures were chosen because they are widely used in the financial market according to Bacon [2008] and Chan [2021].

The annualized return serves as a metric for gauging the average rate of return per year during a designated investment period. Typically computed using the formula:

$$\text{An. Ret.} = \left(\prod_{t=1}^{252} 1 + R_t \right)^{\frac{1}{252}} - 1$$

where R_t denotes the daily return, this measure encapsulates the compounded effect of daily returns over the entire investment horizon. For logarithmic returns r_t , the formula takes a slightly different form:

$$\text{An. Ret.} = \left(\frac{1}{T} \sum_{t=1}^T r_t \right) \times 252$$

This representation for logarithmic returns provides an alternative perspective on the annualized return, particularly when dealing with continuously compounded returns.

This metric quantifies the volatility or risk of the investment by measuring the dispersion of returns around the mean over a one-year period. The annualized Standard Deviation is the

rescaled daily Standard Deviation and can be calculated as

$$\text{An. Std.-Dev.} = \sqrt{\frac{1}{T} \sum_{t=1}^T (R_t - \mu_R)^2} \times \sqrt{252}$$

with μ_R being the mean of returns.

The Sharpe ratio evaluates the risk-adjusted return by comparing the portfolio's excess return over the risk-free rate to its standard deviation. It is a widely used metric in finance to assess the risk-adjusted performance of an investment or portfolio. It quantifies the excess return generated per unit of risk taken. The formula for the Sharpe ratio is given by:

$$\text{SR} = \frac{\mathbb{E}[R] - R_f}{\sigma} \times \sqrt{252}$$

where R_f is the return of the risk-free asset and $\sigma = \sqrt{\frac{1}{T} \sum_{t=1}^T (R_t - \mu_R)^2}$, with μ_R being the mean of returns.

The Sortino ratio is a modified version of the Sharpe ratio that considers only downside risk, which is calculated using the semi-deviation. It tries to provide a more relevant measure of risk-adjusted performance, especially in instances when investors are concerned about downside volatility, represented by $\sigma^- = \sqrt{\frac{1}{T} \sum_{t=1}^T (\min(R_t - \mu_R, 0))^2}$.

$$\text{Sortino} = \frac{\mathbb{E}[R] - R_f}{\sigma^-} \times \sqrt{252}$$

VaR, denoted as $VaR_\alpha(X)$, is a critical risk management metric that assesses the maximum expected loss with a $1 - \alpha$ confidence level over a defined time horizon. It serves as a powerful tool for quantifying the potential downside risk inherent in an investment portfolio.

The calculation for VaR_α involves determining the loss level at which there is a probability of exceeding it during the specified time frame. Mathematically, it can be expressed as follows:

$$VaR_\alpha = \inf\{x \in \mathbb{R} \mid F_X(x) \geq \alpha\}$$

$$\text{An. } VaR_\alpha = VaR_\alpha \times \sqrt{252}.$$

Here, we used VaR at 5% level.

CVaR, denoted as $CVaR_\alpha(X)$, represents a critical risk measure used to gauge the expected loss under extreme scenarios. It is calculated as the expected value of the random variable X conditional on X being less than or equal to its VaR at the same confidence level. This can be

expressed mathematically as follows:

$$\begin{aligned} CVaR_\alpha &= \mathbb{E}[X|X \geq VaR_\alpha(X)] \\ &= \frac{1}{1-\alpha} \int_\alpha^1 VaR_u(X) du \end{aligned}$$

$$An. CVaR_\alpha = CVaR_\alpha \times \sqrt{252}.$$

In essence, $CVaR_\alpha(X)$ provides valuable insights into the potential loss magnitude beyond the VaR, by averaging over all levels for $u \leq \alpha$. The level of CVaR used is 5%.

The worst-drawdown is a measure that quantifies the maximum percentage decline in a portfolio's value from a previous peak to the lowest subsequent point. It helps investors understand the largest loss they might have experienced during a specific investment period. The worst-drawdown can be calculated using the following formula:

$$\text{Worst-Drawdown} = \text{Maximum Drawdown} = \max_{i,j} \left(\frac{V_i - V_j}{V_i} \right) \times 100\%,$$

where V_i is the portfolio's value at time i , V_j is the lowest subsequent value after the peak at time j , and $\max_{i,j}$ represents the maximum value over all peak-to-trough periods.

4 Empirical Analysis

To assess the effectiveness of the proposed strategy, I conducted a rigorous backtest, taking into account statistical biases. Initially, the data was collected and appropriately handled. To avoid model overfitting, cointegration analysis was only performed in the in-sample period, while trading operations were only undertaken in the out-of-sample period.

The assets used were randomly selected from the composition of the Bovespa index, mitigating potential issues related to data-snooping bias and survivorship bias. The generated signals were verified, and entries were made the following day to avoid look-ahead bias, ensuring the portfolio's return was calculated without incorporating future information.

Subsequently, risk and return measures were examined, and the portfolio's performance was compared to a benchmark, the Bovespa index, demonstrating the strategy's success. Transaction costs of 0.03% based on [Frazzini et al. \[2018\]](#) have been included. These steps ensured a thorough examination, considering numerous circumstances that could impact the backtest results.

4.1 Data

The historical data of adjusted close prices for stocks and the composition of the Bovespa index were collected from the Economatica database, covering the period from January 2, 1997, to December 28, 2023. The chosen time frame represents the longest duration obtained since the inception of the Real Plan implementation, an event that marks a significant change in Brazil's monetary and fiscal policies. A total of 203 assets were collected during this period, including the tickers listed in Table 1.

ABEV3	ACES4	AEDU3	ALLL11	ALLL3	ALPA4	AMBV4	AMER3	ARCE3
ARCE4	ARCZ6	ASAI3	ATMP3	AZUL4	B3SA3	BBAS3	BBAS4	BBDC3
BBDC4	BBSE3	BEEF3	BESP4	BHIA3	BIDI11	BISA3	BMTO4	BNCA3
BPAC11	BPAN4	BRAP4	BRDT4	BRFS3	BRKM5	BRML3	BRPR3	BRTP3
BRTP4	CASH3	CCRO3	CESP5	CESP6	CEVA4	CGAS5	CIEL3	CLSC4
CMET4	CMIG3	CMIG4	CMIN3	COGN3	CPFE3	CPLE6	CPSL3	CRFB3
CRTP5	CRUZ3	CSAN3	CSNA3	CSTB4	CTAX4	CTIP3	CVCB3	CYRE3
DASA3	DURA4	DXCO3	EBEN4	EBTP3	EBTP4	ECOR3	EGIE3	ELET3
ELET6	ELPL4	EMAE4	EMBR3	EMBR4	ENBR3	ENEV3	ENGI11	EPTE4
EQTL3	ERIC4	EVEN3	EZTC3	FIBR3	FLRY3	GEPA4	GETI4	GFS3
GGBR4	GNDI3	GOAU4	GOLL4	HAPV3	HGTX3	HYPE3	IGTA3	IGTI11
INEP4	IRBR3	ITSA4	ITUB4	JBSS3	JHSF3	KLBN11	KLBN4	LAME4
LAND3	LCAM3	LIGT3	LIPR3	LOGG3	LREN3	LWSA3	MGLU3	MMXM3
MRFG3	MRVE3	MULT3	NETC4	NTCO3	OGXP3	OIBR3	OIBR4	PALF3
PCAR3	PCAR4	PDGR3	PETR3	PETR4	PETZ3	PMAM4	POMO4	POSI3
PRI03	PRML3	PRTX3	PTIP4	QUAL3	RADL3	RAIL3	RCTB31	RCTB41
RDCD3	RDOR3	RENT3	REPA4	RLOG3	RRRP3	RSID3	RUMO3	SANB11
SAPR11	SBSP3	SDIA4	SHAP4	SLCE3	SMLS3	SOMA3	SUBA3	SULA11
SUZB3	SUZB5	SYNE3	TAEE11	TAMM4	TBLE6	TCOC4	TCSL4	TELB3
TELB4	TIMS3	TLCP4	TMAR5	TMAR6	TMCP4	TNEP4	TNLP3	TNLP4
TOTS3	TPRC6	TRJC6	TRPL4	TSPC3	TSPC6	UBBR11	UGPA3	UGPA4
UNIP6	USIM3	USIM5	VALE3	VALE5	VBBR3	VCPA4	VIVO4	VIVT3
VIVT4	VVAR11	WEGE3	WHMT3	YDUQ3				

Table 1: List of all tickers used to search for cointegrated assets from 1997 to 2023.

The assets were separated between in-sample periods of one year and out-of-sample periods of six months. The dataset is updated every six months. The assets in the in-sample dataset are always chosen because they were part of the Bovespa index composition at the time. The same assets are kept during the out-of-sample period, whether or not they were delisted. Figure 4.1 illustrates this division over time, with the in-sample period in blue and the out-of-sample period in green.

During the backtest execution, priority was given to ensuring that the portfolio contained 10 assets to avoid excessive diversification, as the total amount of assets depends on the investor. The assets were selected randomly. Adhering to the concept that our portfolio is cointegrated, stocks were chosen, and cointegration was tested. If no cointegration relationship was found, an additional 10 assets were randomly selected until the hypothesis was satisfied. If, during the

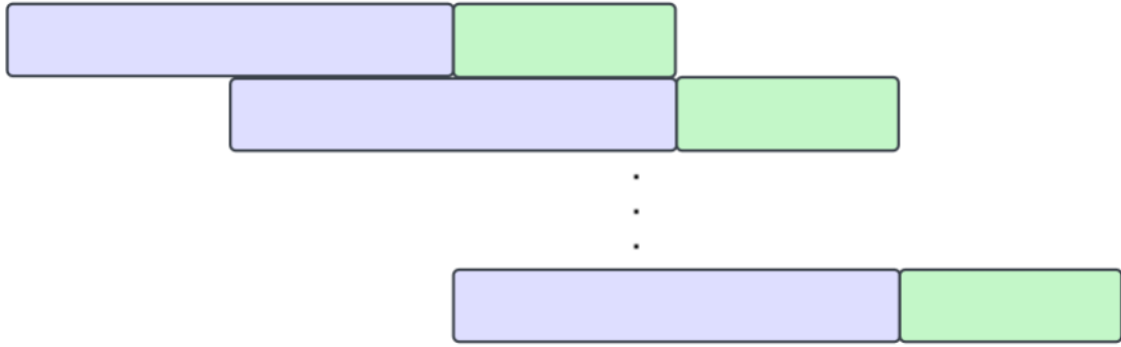


Figure 2: Illustration of period splitting for in-sample and out-of-sample data, with a one-year and six-month rolling window updated every six months.

period, the hypothesis was rejected for all possible combinations, the portfolio’s return for that period would be equal to zero, as no operations would be conducted.

4.2 Results

The empirical analysis compared the results generated by the cointegrated portfolio against a benchmark, the Bovespa index. It was found that the statistical arbitrage strategy achieved superior performance during the analyzed period. Table 2 presents the results for risk and return measures.

Statistics	Portfolio	Benchmark
Annualized Returns	13.9200	10.0900
Annualized Standard Deviation	23.6500	30.4700
Sharpe Ratio	0.5900	0.3300
Sortino Ratio	0.2400	0.4500
Annualized VaR 95%	-0.0000	-0.4600
Annualized CVaR 95%	-0.4100	-0.6900
Worst Drawdown	-0.2300	-0.8000

Table 2: Performance measures comparing cointegrated portfolio against Bovespa index for the period 1998-2023.

According to the reported results, the portfolio’s annualized return is higher than the benchmark. One likely explanation for this performance is the strategy’s capacity to create profits regardless of the economic condition, which leads to increased stability over time. This tendency is readily visible in figure 3, where the portfolio’s development (blue line) contrasts with the benchmark’s extensive period of lateralization (green line). The red and orange lines represent the portfolio’s and benchmark’s drawdowns, respectively.

Furthermore, the cointegrated portfolio exhibits lower volatility, with downturns being found to be less severe than the benchmark. This discovery implies that the portfolio has a faster

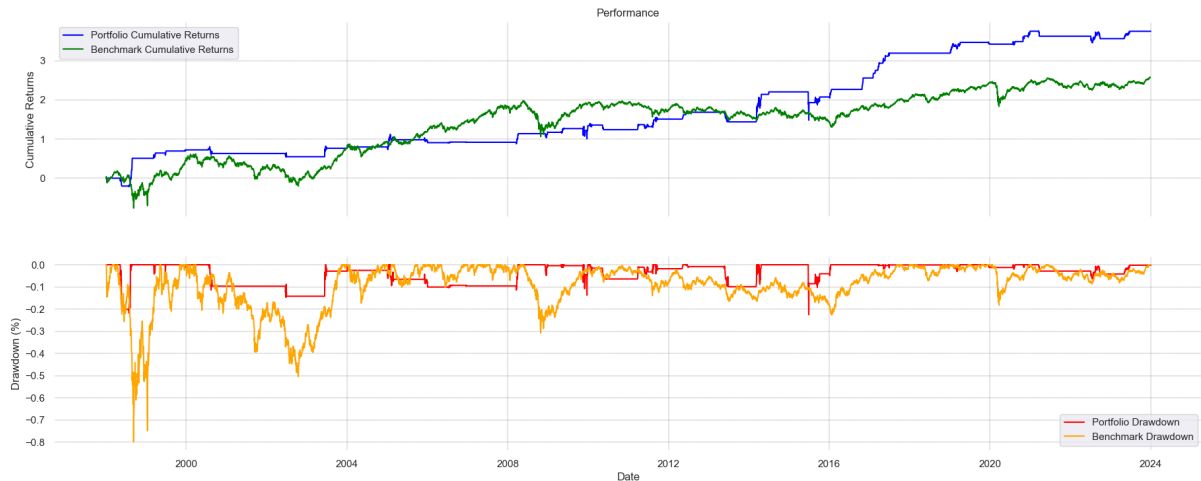


Figure 3: Performance comparison between cointegrated portfolio and Bovespa index.

recovery potential, which is critical for preserving consistency in results. These behavioral patterns demonstrate the strategy’s robustness and efficiency.

The Sharpe and Sortino ratios contribute to this analysis since they reflect the link between expected return and risk, emphasizing the cointegrated portfolio’s superiority. What is noteworthy here is the fact that the portfolio outperformed the Sharpe ratio by a large vantage. Also, the observation that the Sortino ratio does not outperform the Sharpe ratio is a result of higher volatility during downturns - a negative factor - is higher than volatility overall, thereby indicating the risk associated with structural breaks in the spread.

Moreover, while our yearly VaR and CVaR values are lower than the benchmark, our cointegrated portfolio has never experienced decreases at the CVaR level. On the contrary, the Bovespa index fell by more than the value of the CVaR, indicating that the portfolio’s performance is unaffected by market changes.

The findings align with similar studies conducted in the Brazilian market, such as those by [Perlin \[2009\]](#), [Caldeira and Moura \[2013\]](#), and more recently, [Venturini and de Moraes \[2024\]](#). However, the distinction lies in the fact that these studies traditionally focused on comparing the financial arbitrage strategy using pairs of cointegrated assets rather than a comprehensively optimized investment portfolio. This highlights that, by extending the method to the multivariate case, it continues to demonstrate its effectiveness.

4.3 Robustness Test

A critical feature of a backtest is evaluating the outcomes acquired from factors, particularly when working with a statistical arbitrage approach. In this environment, portfolio returns must be independent of the causes underlying market behavior.

To execute a robustness test, portfolio returns must be regressed against these factors. In this regard, I will use the Fama-French five-factor model, as indicated by [Fama and French \[2015\]](#). The dependent variables in this are RMRF (market return relative to the risk-free rate), SMB (small minus big), HML (high minus low), RMW (robust minus weak), and CMA (conservative minus aggressive).

The coefficients associated with each factor are estimated using the Ordinary Least Squares (OLS) method. Given the objective of the test, OLS is well-suited for this purpose due to its simplicity and efficiency in estimating linear relationships between variables, aligning with the scope and requirements of the analysis. Another study employing a similar test is conducted by [Sabino da Silva et al. \[2023\]](#). The regression equation is expressed as follows:

$$R_{port} = -0.0105 \times RMRF + 0.0146 \times SMB - 0.0044 \times HML + 0.0098 \times RMW - 0.0014 \times CMA.$$

Table 3 presents the results of the relevant tests conducted in our analysis.

Test	Test Statistic	P-Value	Conclusion
R-squared	0.000	-	Low variability explanation
Jarque-Bera	11062092.095	0.00	Non-normality of residuals
Omnibus	5411.615	0.00	Low significance between parameters
Durbin-Watson	2.210	-	Possible autocorrelation in residuals

Table 3: Robustness Test Results conducted with Brazilian market factors from 1998 to 2023.

The analysis yields a R^2 equal to zero, indicating that the factors evaluated cannot adequately explain the portfolio's performance. This low explained variability shows that the portfolio approach is genuinely market-neutral. In the meantime, the Omnibus test shows that the explained variance is less than the unexplained variance, which indicates that the model parameters have minimal relevance, validating the notion that the included factors have little influence on the portfolio's performance. Finally, the Jarque-Bera test validates the residuals' non-normality. This lack of normality underscores the presence of patterns or behaviors not captured by the factors, emphasizing additional complexity in the portfolio returns.

5 Conclusion

The decision to use multivariate cointegration to generate a portfolio, rather than just pairs, was based on Johansen's theoretical framework. This methodology enabled the execution of cointegration tests on several assets simultaneously, supporting portfolio design based on the logic of mean-reversion trading through the spread, which represents price divergences from an equilibrium point.

The backtest findings show that statistical arbitrage has the potential to earn rewards while posing no substantial risks. During the investigated time, the strategy produced great returns, outperforming the benchmark, while having lower drawdowns. The robustness test validated the independence of returns from market fluctuations, which is an expected characteristic of an arbitrage strategy.

A few words of caution are also in order. To begin, transaction costs were included in the backtest in addition to the steps taken to avoid statistical biases. Given the portfolio's short positions, these expenses can occasionally result in considerable losses, potentially negating the acquired return. My analysis revealed that even in this case, the cointegrated portfolio would outperform the market.

The pre-selection of assets is another unexplored element that could improve the results. Approaches such as [Sarmiento and Horta \[2021\]](#), which uses unsupervised machine learning algorithms to discover stocks in clusters, are intriguing because they restrict the search, making it easier to locate assets with more common behaviors. It is vital to note that structural breaks in the out-of-sample period may have an impact on the spread, as previously mentioned. Hence, improved spread modeling and forecasting can provide a clearer description of entry and exit points. [Meucci \[2009\]](#) indicates that an Ornstein-Uhlenbeck process properly represents residuals, allowing for a more precise calculation of its mean and variance to produce a Z-score that better reflects price divergences. Other works, such as that of [Robert J. Elliott and Malcolm \[2005\]](#) used the Markov Chain to model the spread, or by [Dunis et al. \[2006\]](#) who employed neural networks to forecast the spread, can be considered to improve the strategy.

All these suggestions can be incorporated into the same framework, providing a more sophisticated and concrete approach. Lastly, analyzing specific sub-periods could add more reliability to the presented results. Given the extensive database, assessing performance during significant events like the dot-com bubble, the 2008 financial crisis, and the COVID-19 pandemic would be interesting.

References

- C.R. Bacon. *Practical Portfolio Performance Measurement and Attribution*. The Wiley Finance Series. Wiley, 2008. ISBN 9780470778050.
- Rodrigo De Losso da Silveira Bueno. *Econometria de séries temporais*. Cengage Learning, 2018.
- João Frois Caldeira and Guilherme Valle Moura. Seleção de uma carteira de pares de ações

- usando cointegração: Uma estratégia de arbitragem estatística. *Brazilian Review of Finance*, 11(1):49–80, 2013.
- Rodolfo C. Cavalcante, Rodrigo C. Brasileiro, Victor L.F. Souza, Jarley P. Nobrega, and Adriano L.I. Oliveira. Computational intelligence and financial markets: A survey and future directions. *Expert Systems with Applications*, 55:194–211, 2016.
- Ernest P Chan. *Quantitative trading: how to build your own algorithmic trading business*. John Wiley & Sons, 2021.
- Richard V Diamond. Learning and trusting cointegration in statistical arbitrage. *Available at SSRN 2220092*, 2014.
- Christian Dunis, Jason Laws, and Ben Evans. Modelling and trading the gasoline crack spread: A non-linear story. *Journal of Derivatives & Hedge Funds*, 12:126–145, 2006.
- Robert F. Engle and C. W. J. Granger. Co-integration and error correction: Representation, estimation, and testing. *Econometrica*, 55(2):251–276, 1987a.
- Robert F. Engle and C. W. J. Granger. Co-integration and error correction: Representation, estimation, and testing. *Econometrica*, 55(2):251–276, 1987b.
- Eugene F. Fama. Efficient capital markets: A review of theory and empirical work. *The Journal of Finance*, 25(2):383–417, 1970.
- Eugene F. Fama and Kenneth R. French. A five-factor asset pricing model. *Journal of Financial Economics*, 116(1):1–22, 2015.
- Andrea Frazzini, Ronen Israel, and Tobias J Moskowitz. Trading costs. *Available at SSRN 3229719*, 2018.
- Evan Gatev, William N. Goetzmann, and K. Geert Rouwenhorst. Pairs Trading: Performance of a Relative-Value Arbitrage Rule. *The Review of Financial Studies*, 19(3):797–827, 02 2006.
- Jesus Gonzalo and Clive Granger. Estimation of common long-memory components in cointegrated systems. *Journal of Business Economic Statistics*, 13(1):27–35, 1995.
- Søren Johansen. Statistical analysis of cointegration vectors. *Journal of Economic Dynamics and Control*, 12(2):231–254, 1988.
- Søren Johansen. Estimation and hypothesis testing of cointegration vectors in gaussian vector autoregressive models. *Econometrica*, 59(6):1551–1580, 1991.

- Andrew W. Lo and A.Craig MacKinlay. The size and power of the variance ratio test in finite samples: A monte carlo investigation. *Journal of Econometrics*, 40(2):203–238, 1989.
- Gangadharrao S Maddala and In-Moo Kim. Unit roots, cointegration, and structural change. 1998.
- Attilio Meucci. Review of statistical arbitrage, cointegration, and multivariate ornstein-uhlenbeck. 2009.
- Daniel P Palomar. Pair’s trading. *Portfolio Optimization with R*, 2020.
- Marcelo Scherer Perlin. Evaluation of pairs-trading strategy at the brazilian financial market. *Journal of Derivatives & Hedge Funds*, 15:122–136, 2009.
- J.P. Ramos-Requena, J.E. Trinidad-Segovia, and M.A. Sánchez-Granero. Introducing hurst exponent in pair trading. *Physica A: Statistical Mechanics and its Applications*, 488:39–45, 2017.
- John Van Der Hoek * Robert J. Elliott and William P. Malcolm. Pairs trading. *Quantitative Finance*, 5(3):271–276, 2005.
- Fernando A.B. Sabino da Silva, Flavio A. Ziegelmann, and João F. Caldeira. A pairs trading strategy based on mixed copulas. *The Quarterly Review of Economics and Finance*, 87:16–34, 2023.
- Simão Moraes Sarmiento and Nuno Horta. *A Machine Learning based Pairs Trading Investment Strategy*. Springer, 2021.
- Stephen J Taylor. *Asset price dynamics, volatility, and prediction*. Princeton university press, 2011.
- Ariel Amadeu Edwards Teixeira. *Pair trading in Bovespa with a quantitative approach: cointegration, Ornstein-Uhlenbeck equation and Kelly criterion*. PhD thesis, 2014.
- Lucas Stumpf Venturini and Gustavo Inácio de Moraes. Pairs trading aplicado ao mercado de capitais brasileiro: Uma abordagem via cointegração e ornstein-uhlenbeck. *Revista de Administração, Contabilidade e Economia da Fundace*, 15(1), 2024.
- Wenjun Xie, Zhao Zhi Toh, and Yuan Wu. Copula-based pairs trading in asia-pacific markets. , 24(4):1–17, 2016.