



CERTIFICATE IN QUANTITATIVE FINANCE  
FINAL PROJECT

**MULTIVARIATE COINTEGRATION IN STATISTICAL  
ARBITRAGE**  
JOÃO RAMOS JUNGBLUT

**2024**



# Contents

1	INTRODUCTION . . . . .	1
2	METHODOLOGY . . . . .	4
2.1	Model Specification . . . . .	4
2.1.1	Vector Autoregressive (VAR) . . . . .	4
2.1.2	Cointegration . . . . .	5
2.1.3	Vector Error Correction Model (VECM) . . . . .	6
2.1.4	Johansen's approach . . . . .	9
2.1.5	Maximum Likelihood Estimation of $\Phi$ . . . . .	10
2.2	Trading Strategy . . . . .	13
2.3	Performance Measures . . . . .	14
3	EMPIRICAL ANALYSIS . . . . .	18
3.1	Data . . . . .	18
3.2	Results . . . . .	20
3.3	Robustness Test . . . . .	21
4	CONCLUSION . . . . .	23
	REFERENCES . . . . .	25
	APPENDIX - PYTHON CODE . . . . .	28



# 1 INTRODUCTION

Arbitrage 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.

The underlying assumption behind such a concept is that financial markets are efficient, as per Fama [1970]. The hypothesis suggests that markets reflect all available information, making profit opportunities short-lived because participants tend to act quickly to exploit these differences, restoring price equilibrium. Therefore, it is easier to find arbitrage by trading multiple assets, precisely because it is more challenging to model hidden patterns in multiple variables than in just one.

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. This means that short-term linear relationships, characterized by metrics such as correlation, are inconsistent and unreliable. Therefore, 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 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 relationship, this series tends to revert to the mean in the long term, allowing for long and short oper-

ations. 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 [1988], 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, what I will do in this project is build a cointegrated portfolio using Johansen [1988] test which allows for exploring cointegration in a multivariate case.

I will use stock price data comprising the Bovespa index from 1996 to 2023, collected from Economática. The database will be consistently divided into one year in-sample periods and six month out-of-sample periods. I will seek cointegrated eigenvectors using the

Johansen test and construct stationary residuals with the aim of profiting from statistical arbitrage operations.

The structure of this work is organized as follows: in Chapter 2, I detail in Section 2.1 the specifications of VAR, Cointegration, VECM, Johansen Approach, and Maximum Likelihood Estimation. In Section 2.2, the strategy is presented, including the definition of entry and exit points, as well as the calculation of the percentage of each asset in the portfolio. Section 2.3 discusses the measures used to evaluate the portfolio's performance. In the following chapter, 3, I present the data used in 3.1 more clearly and the results in Section 3.2, along with the robustness test in 3.3. Finally, in the conclusion, Chapter 4, I revisit the results and fundamental concepts of the model to discuss improvements and future work. The codes are available in Appendix.

## 2 METHODOLOGY

### 2.1 Model Specification

The models for implementing the trading strategy will be presented here. First, I will explain the Vector Autoregressive (VAR), which serves as the basis for a deeper understanding of the others. Following that, I will address the concept of cointegration, and subsequently formulate the Vector Error Correction Model (VECM). The Johansen test will be used to identify cointegration vectors, and parameter estimation will then be performed through Maximum Likelihood. Finally, I will highlight how the strategy is executed through the spread, and the risk-return measures used to evaluate portfolio performance. The notations used are based on the works of Bueno [2018].

#### 2.1.1 Vector Autoregressive (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  $A$ :

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

$A$  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  $\varepsilon_t$  is a  $n \times 1$  vector of uncorrelated random disturbances.

Examining the explanation for a bivariate situation helps to better understand the



model's endogeneity. Consider these equations:

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

In this context, (i)  $x_{1,t}$  and  $x_{2,t}$  are stationary, (ii)  $\varepsilon_{x_1,t} \sim RB(0, 1)$  and  $\varepsilon_{x_2,t} \sim RB(0, 1)$ , and (iii)  $\varepsilon_{x_1,t} \perp \varepsilon_{x_2,t} \Rightarrow Cov(\varepsilon_{x_1,t}, \varepsilon_{x_2,t}) = 0$ .

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:

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

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  $Ae_t = B\varepsilon_t$ . 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.

### 2.1.2 Cointegration

According to Engle and Granger [1987], 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  $\beta$  such that

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

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

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

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

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

With the normalization of  $\beta_2 = \frac{\tilde{\beta}_2}{\tilde{\beta}_1}$ ,  $u_t$  can be considered the residual of one coordinate of  $X_t$  against the other variables. Thus, the variables are cointegrated, and  $\beta$  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.

An important observation is that the deterministic trend dominates first-order non-stationarity, which, in turn, dominates stationarity in the sense of asymptotic convergence.

### 2.1.3 Vector Error Correction Model (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:

$$\begin{aligned}\Delta x_{1,t} &= \alpha_1 \hat{u}_{t-1} + \sum_{j=1}^{p-1} \lambda_{11,j+1} \Delta x_{1,t-j} + \sum_{j=1}^{p-1} \lambda_{12,j+1} \Delta x_{2,t-j} + e_{x_1,t} \\ \Delta x_{2,t} &= \alpha_2 \hat{u}_{t-1} + \sum_{j=1}^{p-1} \lambda_{21,j+1} \Delta x_{1,t-j} + \sum_{j=1}^{p-1} \lambda_{22,j+1} \Delta x_{2,t-j} + e_{x_2,t}.\end{aligned}$$

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

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

A 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. This allows us to state Granger's theorem.

**Theorem 1 (Granger's Theorem)** *If  $|\Phi(Z)| = 0$  implies  $Z > I$  and  $0 < \text{rank}(\Phi) = r < n$ , then there exist matrices  $\alpha$  and  $\beta$  with dimensions  $n \times r$  such that:  $\Phi = \alpha\beta$ .*

- $\beta$  is the matrix of **cointegration**.
- $\alpha$  is the matrix of **adjustment**.

The theorem expresses the idea that  $\Phi$  can be decomposed into two multiplicative matrices. From this, we derive the VECM, recursively adding and subtracting past terms to generalize the equation. Thus, let

$$\begin{aligned}
X_t &= \Phi_1 X_{t-1} + \Phi_2 X_{t-2} + \Phi_3 X_{t-3} + e_t \\
&= \Phi_1 X_{t-1} + \Phi_2 X_{t-2} + \boxed{\Phi_3 X_{t-2} - \Phi_3 X_{t-2}} + \Phi_3 X_{t-3} + e_t \\
&= \Phi_1 X_{t-1} + (\Phi_2 + \Phi_3) X_{t-2} - \Phi_3 \Delta X_{t-2} + e_t \\
&= \Phi_1 X_{t-1} + \boxed{(\Phi_2 + \Phi_3) X_{t-1} - (\Phi_2 + \Phi_3) X_{t-1}} + (\Phi_2 + \Phi_3) X_{t-2} - \Phi_3 \Delta X_{t-2} + e_t \\
&= (\Phi_1 + \Phi_2 + \Phi_3) X_{t-1} + (\Phi_2 + \Phi_3) \Delta X_{t-1} - \Phi_3 \Delta X_{t-2} + e_t \\
X_t \boxed{-X_{t-1}} &= \boxed{-X_{t-1}} + (\Phi_1 + \Phi_2 + \Phi_3) X_{t-1} + (\Phi_2 + \Phi_3) \Delta X_{t-1} - \Phi_3 \Delta X_{t-2} + e_t \\
\Delta X_t &= -[I - (\Phi_1 + \Phi_2 + \Phi_3)] X_{t-1} + (\Phi_2 + \Phi_3) \Delta X_{t-1} - \Phi_3 \Delta X_{t-2} + e_t \\
&= \Phi X_{t-1} + \sum_{i=1}^2 \Lambda_i \Delta X_{t-i} + e_t,
\end{aligned}$$

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

The model is called *error correction* because it explains  $\Delta X_t$  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.

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

### 2.1.4 Johansen's approach

Johansen presents a test to identify the rank of the matrix  $\Phi$  and, consequently, estimate the cointegration vectors included in the matrix  $\beta$ . 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  $r$  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  $\alpha$  and  $\beta$  are not identifiable as they form an overparameterization of the model. However, we can delimit the cointegration space to  $\text{span}(\beta)$ .

### 2.1.5 Maximum Likelihood Estimation of $\Phi$

The estimation of  $\Phi$  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  $\Pi_{ij} = \frac{\sum_{t=1}^T \Upsilon_{i,t}\Upsilon'_{j,t}}{T}$ , we can rewrite the previous equation as:

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

I will show that to obtain  $e_t$ , we do not need  $\hat{\Lambda}$ .

First, let's regress  $\Upsilon_{0,t}$  on  $\Upsilon_{2,t}$ , that is,

$$\Upsilon_{0,t} = B\Upsilon_{1,t} + r_{0,t}$$

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, I will regress  $\Upsilon_{1,t}$  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:

$$\begin{aligned} \hat{e}_t &= \Upsilon_{0,t} - \alpha\beta'r_{1,t} - (\Pi_{02}\Pi_{22}^{-1} - \alpha\beta'\Pi_{12}\Pi_{22}^{-1})\Upsilon_{2,t} \\ &= \Upsilon_{0,t} - \Pi_{02}\Pi_{22}^{-1}\Upsilon_{2,t} - \alpha\beta'(\Upsilon_{1,t} - \Pi_{12}\Pi_{22}^{-1}\Upsilon_{2,t}) \\ &= r_{\hat{0},t} - \alpha\beta'r_{\hat{1},t}. \end{aligned}$$

Thus, the quantities  $r_{\hat{0},t}$  and  $r_{\hat{1},t}$  can be obtained from the auxiliary regressions, and we can estimate  $\alpha$  and  $\beta$ . The new function to maximize becomes:

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

Equivalent to maximizing this function, one can perform OLS regression between  $r_{\hat{0},t}$  and

$r_{1,t}^\wedge$ . One way is to estimate  $\alpha$  and  $\Sigma$  for a given  $\beta$ ,

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

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

with  $S = T^{-1} \sum_{t=1}^T r_{i,t}^\wedge r_{j,t}^\wedge$ ,  $j = 0, 1$ .

The matrix  $\beta$  has not been estimated yet; however, starting from the last derivation,  $\beta$  can be found by maximizing the likelihood.

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

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

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

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  $\alpha$  and  $\beta$  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  $\Delta X_t$  (trend). Furthermore, as we are only looking at the cointegrated space, the total residual presented will be

$$e_t = \beta'X_t.$$



## 2.2 Trading Strategy

The chosen strategy is similar to pairs trading using the cointegration method. Our approach, on the other hand, 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 an equilibrium.

Detecting these price divergences requires observing the residuals of the linear relationships. Although the strategy is almost identical to pairs trading, in the multivariate model, we have  $n$  financial assets. As explained in Section 2.1, the residuals are given by:

$$e_t = \beta' X_t.$$

To clarify the strategy's implementation, we shall use a separate nomenclature here. During the out-of-sample phase, the computed residuals will be referred to as "spread" and will 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  $Z_{\text{score}} = \frac{\text{Spread} - \mu}{\sigma}$ .

Meanwhile, the normalized cointegrated vector  $\beta = [1, \frac{\tilde{\beta}_2}{\tilde{\beta}_1}, \frac{\tilde{\beta}_3}{\tilde{\beta}_1}, \dots, \frac{\tilde{\beta}_n}{\tilde{\beta}_1}]$  will constitute the hedge ratios. To determine the portfolio weights  $W$ , we divide the hedge ratios by their norm.

$$W = \frac{\beta}{\|\beta\|}$$

This way, we secure two critical points: first, that only available capital is used; second, that overly significant percentages are not concentrated in a single asset. The portfolio return is calculated as:

$$R_{\text{port}} = WR_t.$$

The strategy is built on the spread 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 in this context because it allows for the visualization of mean-reverting movements and determines entry and exit locations based on standard deviations.

I have used the following parameters:

- **Entry Point:** 1 standard deviation;
- **Exit Point:** 0 (mean);
- **Stop Loss:** 3 standard deviations.

When the Z-score reaches -1, we open a long position in the portfolio, closing it with a positive return at a Z-score of 0 or with a loss at -3. Conversely, as the Z-score achieve 1, we open a short position. The figure 1 depicts these locations, with the dashed green lines representing entrance points, the dashed black lines representing exit points, the dashed red line representing the Stop Loss, and the blue line representing the spread.

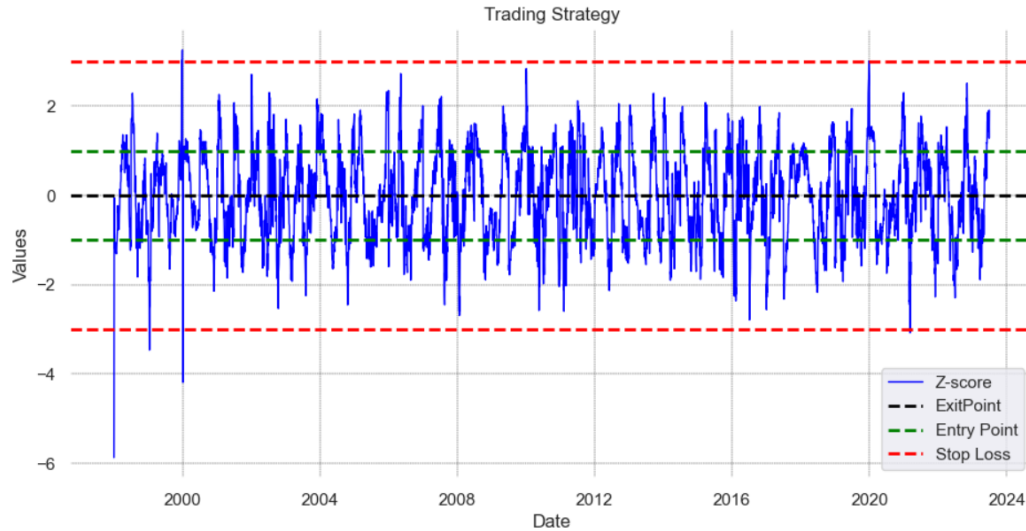


Figure 1: Trading Strategy

## 2.3 Performance Measures

The evaluation of the investment portfolio's performance involves the computation of diverse measures. These metrics offer a thorough evaluation of 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, Sortino ratio,  $VaR_{0.95}$ ,  $CVaR_{0.95}$  and Worst-Drawdown.

## Annualized Return

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:

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

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

## Annualized Standard Deviation

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

$$\text{An. Std.-Dev.} = \sigma \times \sqrt{252}.$$

## Sharpe Ratio

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:

$$SR = \frac{E[R] - R_f}{\sigma}.$$

## Sortino Ratio

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.

$$\text{Sortino} = \frac{E[R] - R_f}{\sigma^-}.$$

## Value at Risk (VaR)

Value at Risk, denoted as  $VaR_\beta(X)$ , is a critical risk management metric that assesses the maximum expected loss with a  $1 - \beta$  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  $VaR_\beta$  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:

$$\begin{aligned} Pr(x \leq VaR(X)) &= 1 - \beta \\ VaR_\beta(X) &= \inf\{x | Pr(X > x) \leq 1 - \beta\} \\ &= \inf\{x | F_X(x) \geq \beta\} \end{aligned}$$

Here, I will be using VaR at 5% level.

## Conditional Value at Risk (CVaR)

Conditional Value at Risk, denoted as  $CVaR_\beta(X)$ , 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 Value at Risk ( $VaR$ ) at the same confidence level. This can be expressed mathematically as follows:

$$\begin{aligned} CVaR_\beta(X) &= E[X|X \leq VaR_\beta(X)] \\ &= \frac{1}{1-\beta} \int_0^{1-\beta} VaR_\alpha(X) d\alpha \end{aligned} \quad (1)$$

In essence,  $CVaR_\beta(X)$  provides valuable insights into the potential loss magnitude beyond the  $VaR_\beta(X)$  threshold, taking into account extreme scenarios with a confidence level of  $1 - \beta$ . The level of CVaR used is 5%.

## Worst-Drawdown

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.

### 3 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 on the following day to avoid look-ahead bias, ensuring that 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 were taken to ensure a thorough examination, considering numerous circumstances that could impact the backtest results.

#### 3.1 Data

The historical data of adjusted close prices for stocks and the composition of the Bovespa index were collected from Economática database, covering the period from January 2, 1997, to December 28, 2023. A total of 203 assets were used during this period, including the tickers listed in Table 1.

The assets were separated between in-sample periods of one year and out-of-sample intervals 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 2 illustrates this division over time, with the in-sample period in blue and the out-of-sample period in green.

During the backtest execution, the priority was given to ensuring that the portfolio contained a total of 10 assets to avoid excessive diversification. The assets were selected

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	B RTP3
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	GEP A4	GETI4	GFSA3
GGBR4	GNDI3	GOAU4	GOLL4	HAPV3	HGT X3	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
PRIO3	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	TAE E11	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 Assets

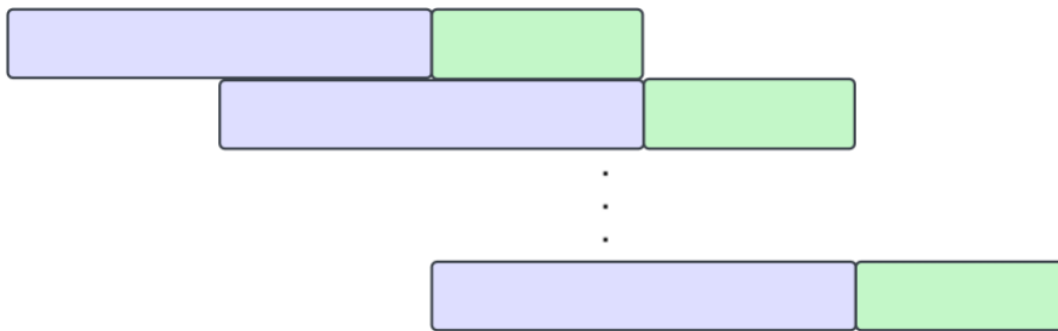


Figure 2: Period Splitting

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 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.

### 3.2 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.

	<b>Portfolio</b>	<b>Benchmark</b>
Annualized Returns	23.1200	10.0900
Annualized Standard Deviation	21.0200	30.4700
Sharpe Ratio	1.1000	0.3300
Sortino Ratio	1.2500	0.4500
Annualized VaR 95%	-0.3100	-0.4600
Annualized CVaR 95%	-0.4900	-0.6900
Worst Drawdown	-0.3000	-0.8000

Table 2: Performance Metrics

According to the reported results, the portfolio's annualized return is much higher than the benchmark. One likely explanation for this exceptional 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 contrasts with the benchmark's extensive period of lateralization.

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



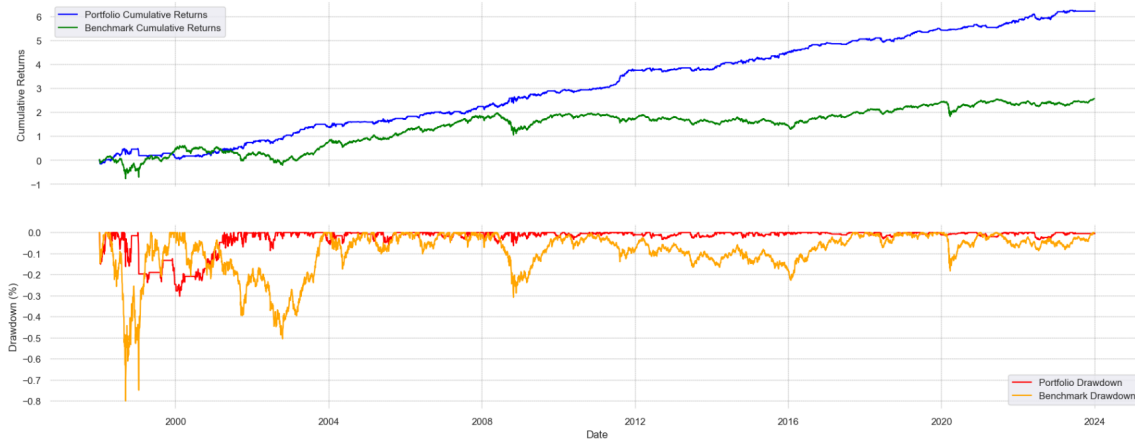


Figure 3: Performance

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 not only the fact that the portfolio outperformed both indices, but also the observation that the Sortino ratio outperforms the Sharpe ratio. As a result, we underline that volatility during downturns, a negative factor, is lower than volatility overall, thus showing risk reduction.

Moreover, while having yearly VaR and CVaR values lower than the benchmark, our cointegrated portfolio has never experienced decreases at these levels. 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 as a whole.

### 3.3 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. The regression equation is expressed as follows:

$$R_{port} = 0.1573 \times RMRF - 0.0519 \times SMB + 0.0400 \times HML - 0.0491 \times RMW + 0.0175 \times CMA.$$

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

Test	Test Statistic	P-Value	Conclusion
R-squared	0.051	-	Low variability explanation
Jarque-Bera	10823.817	0.00	Non-normality of residuals
Omnibus	951.836	0.00	Low significance between parameters
Durbin-Watson	1.979	-	Possible autocorrelation in residuals

Table 3: Robustness Test Results

The analysis yields a  $R^2$  near 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.

## 4 CONCLUSION

The decision to use the statistical arbitrage technique to generate a cointegrated portfolio, rather than just pairs, was based on the Johansen's theoretical framework. This methodology enabled the execution of cointegration tests on several assets at the same time, supporting portfolio design based on the logic of trading through the spread, which represents price divergences from an equilibrium point.

The backtest findings show that 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 a result, improved spread modeling and forecasting can provide a clearer description of entry and exit points, eliminating arbitrary techniques. 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 using the Markov Switch Regime Model to trade in different regimes, or by employing neural networks to predict when the spread returns to the mean, 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, it would be interesting to assess performance during significant events like the dot-com bubble, the 2008 financial crisis, and the COVID-19 pandemic.

## REFERENCES

### References

- Rodrigo de Losso da Silveira Bueno. Econometria de séries temporais. 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 LF Souza, Jarley P Nobrega, and Adriano LI 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.
- Fernando AB 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.
- Richard V Diamond. Learning and trusting cointegration in statistical arbitrage. *Available at SSRN 2220092*, 2014.
- Economática. Economática: Banco de dados. <http://www.economatica.com.br/index fla.htm>.
- Robert F Engle and Clive WJ Granger. Co-integration and error correction: representation, estimation, and testing. *Econometrica: journal of the Econometric Society*, pages 251–276, 1987.
- 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, 2006.
- Søren Johansen. Statistical analysis of cointegration vectors. *Journal of economic dynamics and control*, 12(2-3):231–254, 1988.
- Søren Johansen. Estimation and hypothesis testing of cointegration vectors in gaussian vector autoregressive models. *Econometrica: journal of the Econometric Society*, pages 1551–1580, 1991.
- Andrew W Lo and A Craig MacKinlay. Stock market prices do not follow random walks: Evidence from a simple specification test. *The review of financial studies*, 1(1):41–66, 1988.
- 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.
- José Pedro Ramos-Requena, JE Trinidad-Segovia, and MA Sánchez-Granero. Introducing hurst exponent in pair trading. *Physica A: statistical mechanics and its applications*, 488:39–45, 2017.
- 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.

Wenjun Xie, Zhao Zhi Toh, and Yuan Wu. Copula-based pairs trading in asia-pacific markets. , 24(4):1–17, 2016.

## APPENDIX - PYTHON CODE

### Import Packages

```
1  # Data Manipulation
2  import pandas as pd
3  import numpy as np
4  import random
5
6  # set display options
7  pd.options.display.float_format = "{:,.4f}".format
8  pd.set_option('display.max_columns', 100)
9
10 # Visualizaiton
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 plt.style.use('fivethirtyeight')
14 sns.set(rc={'figure.figsize': (20, 8)})
15
16 # Cointegration test
17 from statsmodels.tsa.vector_ar.vecm import coint_johansen
18
19 # Robust analysis
20 import statsmodels.api as sm
21
22 # Ignore warnings
23 import warnings
24 warnings.filterwarnings('ignore')
```



## Data Manipulation

```
1  # Load stock price data
2  AdjClose = pd.read_excel('StockPrice.xlsx', skiprows=[1],
   ↪   index_col=0, parse_dates=True)
3  AdjClose.index.name = 'date'
4  AdjClose.head()
```

```
1  # Load Ibovespa market composition
2  MarketComposition = pd.read_excel("MarketComposition.xlsx")
3  MarketComposition.head()
```

```
1  # Splitting in-sample and out-of-sample data
2  df = AdjClose.copy()
3
4  inSample = {}
5  outofSample = {}
6
7  update_date = pd.Timestamp('1998-01-01')
8  count = 1
9  while update_date <= pd.Timestamp('2024-01-01'):
10     # In-sample is determined by market composition date
11     inSample_mask = (df.index > (update_date -
   ↪   pd.DateOffset(months=12))) & (df.index <= (update_date))
12     inSample[count] = df.loc[inSample_mask]
13
14     # Out-of-sample
15     outofSample_mask = (df.index > (update_date)) & (df.index <=
   ↪   (update_date + pd.DateOffset(months=6)))
16     outofSample[count] = df.loc[outofSample_mask]
17
```

```

18     # Update rule
19     count += 1
20     update_date += pd.DateOffset(months=6)
21
22     # In-sample tail
23     print("In-sample tail:\n", inSample[1].iloc[:, :5].tail())
24
25     # Out-of-sample head
26     print("Out-of-sample head:\n", outofSample[1].iloc[:, :5].head())

```

```

1     # Keep only stocks from the Bovespa index
2     for index, date in
    → enumerate(MarketComposition['Date'].unique().tolist()):
3         # Select symbols from market composition
4         symbols = MarketComposition.loc[MarketComposition['Date'] ==
    → date, 'Symbols'].tolist()
5         # Filter stocks and clean data
6         inSample[index+1] = inSample[index+1][symbols].ffill()
7         inSample[index+1] = inSample[index+1].dropna(axis=1, how='any')
8         outofSample[index+1] =
    → outofSample[index+1][inSample[index+1].columns].fillna(method='pad')
    → # handle delisted stocks
9
10    # In-sample tail
11    print("In-sample tail:\n", inSample[1].iloc[:, :5].tail())
12
13    # Out-of-sample head
14    print("Out-of-sample head:\n", outofSample[1].iloc[:, :5].head())

```

## Johansen Approach

```
1 class Johansen:
2     """
3     This class performs the Johansen cointegration test and provides
4     → related calculations.
5
6     Attributes:
7     data : array-like
8         Time series data for the cointegration test.
9     det_order : int, optional
10        Order of deterministic terms in the cointegration test.
11    k_ar_diff : int, optional
12        Number of lags to include in the cointegration test.
13    ct : JohansenTestResult, optional
14        Result of the cointegration test.
15    """
16
17    def __init__(self, data, det_order=0, k_ar_diff=1):
18        """
19        Initializes the Johansen class with the provided parameters.
20
21        Parameters:
22        data : array-like
23            Time series data for the cointegration test.
24        det_order : int, optional
25            Order of deterministic terms in the cointegration test.
26        k_ar_diff : int, optional
27            Number of lags to include in the cointegration test.
28        """
29        self.data = data
```

```

29     self.det_order = det_order
30     self.k_ar_diff = k_ar_diff
31     self.ct = coint_johansen(self.data,
    ↪     det_order=self.det_order, k_ar_diff=self.k_ar_diff)
32
33     def n_coint_vectors(self):
34         """
35         Calculates the number of cointegrated vectors based on the
    ↪     cointegration test result.
36
37         Returns:
38         n_vectors : int
39             Number of cointegrated vectors.
40         """
41         return len(np.where(self.ct.lr2 > self.ct.cvm[:, 1])[0])
42
43     def beta(self):
44         """
45         Calculates the cointegrated eigenvectors (beta) based on the
    ↪     cointegration test result.
46
47         Returns:
48         beta : array-like
49             Cointegrated eigenvectors.
50         """
51         return self.ct.evec[:, np.where(self.ct.lr2 > self.ct.cvm[:,
    ↪     1])[0]]
52
53     def alpha(self):
54         """

```

```

55         Calculates the estimated alpha matrix based on the
56         → cointegration test result.
57
58         Returns:
59         alpha : array-like
60         Estimated alpha matrix.
61         """
62         r0t = np.transpose(self.ct.r0t)
63         r1t = np.transpose(self.ct.rkt)
64
65         S01 = np.dot(r0t, r1t.T) / r0t.shape[1]
66         S11 = np.dot(r1t, r1t.T) / r0t.shape[1]
67
68         beta = self.beta()
69         return S01 @ beta @ np.linalg.inv(beta.T @ S11 @ beta)
70
71     def Phi(self):
72         """
73         Calculates the Phi matrix based on the cointegration test
74         → result.
75
76         Returns:
77         Phi : array-like
78         Phi matrix.
79         """
80         return np.dot(self.alpha(), self.beta().T)
81
82     def residuals(self):
83         """
84         Calculates the residuals based on the cointegration test
85         → result.

```

```

83
84     Returns:
85     residuals : array-like
86     Residuals.
87     """
88     return np.dot(self.Phi(), self.data.T)

```

## Backtest Class

```

1  class Backtest:
2      def __init__(self, inSample, outofSample, entry_point = 2,
3          ↪ exit_point=0, stop_loss=None, det_order=0, k_ar_diff=1):
4          """
5          Initializes the Backtest class.
6
7          Parameters:
8          data : DataFrame
9              DataFrame containing the financial data.
10             det_order : int, optional
11                 Order of deterministic components in the cointegration
12                 ↪ test (default is 0).
13             k_ar_diff : int, optional
14                 Number of lags to difference the data in the
15                 ↪ cointegration test (default is 1).
16             """
17             self.inSample = inSample
18             self.outofSample = outofSample
19             self.entry_point = entry_point
20             self.exit_point = exit_point
21             self.stop_loss = stop_loss

```

```

19     self.det_order = det_order
20     self.k_ar_diff = k_ar_diff
21     self.johansen = Johansen(inSample, det_order, k_ar_diff)
22     self._hedge_ratio = self.hedge_ratio()
23
24     def hedge_ratio(self):
25         """
26         Calculates the hedge ratio (Phi) based on the Johansen
27         ↪ cointegration analysis.
28
29         Returns:
30         hedge_ratio : array-like
31             Hedge Ratios.
32         """
33         HR = self.johansen.beta().T
34
35         return HR / HR[0]
36
37     def weight(self):
38         """
39         Calculates the weights for each asset in the portfolio.
40
41         Returns:
42         weights : array-like
43             Portfolio weights.
44         """
45         return np.round(self._hedge_ratio /
46             ↪ np.sum(self._hedge_ratio), 2)[0]
47
48     def spread(self):

```

```

48     """
49     Calculates the spread based on the hedge ratio and the
    ↪ financial data.
50
51     Returns:
52     spread : array-like
53         Spread.
54     """
55     return np.nan_to_num(np.dot(self._hedge_ratio ,
    ↪ self.outofSample.T))[0]
56
57 def z_score(self):
58     """
59     Calculates the z-score of the spread.
60
61     Returns:
62     z_score : DataFrame
63         Z-scores.
64     """
65     mean = np.mean(self.spread())
66     std = np.std(self.spread())
67
68     return (self.spread() - mean) / std
69
70 def signal(self):
71     """
72     Generates trading signals based on the calculated z-score.
73
74     Returns:
75     signals : DataFrame

```



```

76         DataFrame containing the generated trading signals (1:
77         ↪ long, -1: short, 0: no signal).
78         """
79         position = 0
80         signals = []
81         z_score = self.z_score()
82
83         # Loop through the data to generate signals
84         for i in range(0, len(z_score)):
85
86             # Generate signals based on z-score values
87             # Entry points
88             if t <= -self.entry_point and position == 0:
89                 position = 1 # open position
90                 signals.append(position)
91             elif t >= self.entry_point and position == 0:
92                 position = -1 # open position
93                 signals.append(position)
94             # Exit points
95             elif t >= -self.exit_point and (position == 1 or
96             ↪ position == 2):
97                 position = 0 # close position
98                 signals.append(position)
99             elif t <= self.exit_point and (position == -1 or
100             ↪ position == -2):
101                 position = 0 # close position
102                 signals.append(position)
103             # Stop loss
104             elif self.stop_loss is not None:
105                 if t <= -self.stop_loss and position == 1:

```

```

104         position = 2 # close position
105         signals.append(position)
106     elif t >= self.stop_loss and position == -1:
107         position = -2 # close position
108         signals.append(position)
109         # Hold signal
110     else:
111         signals.append(position)
112         # Hold signal
113     else:
114         signals.append(position)
115
116     return pd.DataFrame({"Signal": signals},
117         ↳ index=self.outofSample.index).shift(1)
118
119 def port_rtn(self):
120     """
121     Calculates the daily returns of the portfolio based on the
122     ↳ generated signals.
123
124     Returns:
125     port_returns : DataFrame
126         DataFrame containing the calculated daily portfolio
127         ↳ returns.
128     """
129     stock_rtn = np.log(self.outofSample).diff().fillna(0)
130     port_rtn = []
131     signal = self.signal()
132
133     # Loop through the data to calculate portfolio returns
134     for i in range(0, len(self.outofSample)):

```

```

132         if signal.iloc[i, 0] == 1:
133             rtn = np.dot(self.weight(), stock_rtn.iloc[i]) -
                ↪ 0.0003 # minus transaction costs
134         elif signal.iloc[i, 0] == -1:
135             rtn = - np.dot(self.weight(), stock_rtn.iloc[i]) -
                ↪ 0.0003 # minus transaction costs
136         else:
137             rtn = 0
138
139         port_rtn.append(rtn)
140
141     return pd.DataFrame({"Returns": port_rtn},
        ↪ index=self.outofSample.index)

```

## Performance Measures

```

1 class Performance:
2     def __init__(self, portfolio, benchmark=None):
3         """
4         Initializes the Performance class with portfolio returns and
5         ↪ an optional benchmark.
6
7         Parameters:
8         portfolio : pandas.Series
9             Time series of portfolio returns.
10        benchmark : pandas.Series, optional
11            Time series of benchmark returns (default is None).
12        """
13        self.portfolio = portfolio
14        self.benchmark = benchmark

```

```

14
15 def annual_rtn(self):
16     """
17     Calculates the annualized returns of the portfolio and
18     ↪ benchmark.
19
20     Returns:
21     tuple
22     Tuple containing the annualized portfolio returns and
23     ↪ benchmark returns.
24     """
25     _portfolio = round(100 * (np.mean(self.portfolio) * 252), 2)
26     _benchmark = None
27     if self.benchmark is not None:
28         _benchmark = round(100 * (np.mean(self.benchmark) *
29         ↪ 252), 2)
30     return _portfolio, _benchmark
31
32 def annual_std(self):
33     """
34     Calculates the annualized standard deviation of the
35     ↪ portfolio and benchmark.
36
37     Returns:
38     tuple
39     Tuple containing the annualized portfolio standard
40     ↪ deviation and benchmark standard deviation.
41     """
42     _portfolio = round(100 * (np.std(self.portfolio) *
43     ↪ np.sqrt(252)), 2)[0]
44     _benchmark = None

```

```

39     if self.benchmark is not None:
40         _benchmark = round(100 * (np.std(self.benchmark) *
41             ↪ np.sqrt(252)), 2)[0]
42     return _portfolio, _benchmark
43
44 def sharpe_ratio(self):
45     """
46     Calculates the Sharpe ratio of the portfolio and benchmark.
47
48     Returns:
49     tuple
50         Tuple containing the Sharpe ratio of the portfolio and
51         ↪ benchmark.
52     """
53     _portfolio = round(self.annual_rtn()[0] /
54         ↪ self.annual_std()[0], 2)
55     _benchmark = None
56     if self.benchmark is not None:
57         _benchmark = round(self.annual_rtn()[1] /
58             ↪ self.annual_std()[1], 2)
59     return _portfolio, _benchmark
60
61 def sortino_ratio(self):
62     """
63     Calculates the Sortino ratio of the portfolio and benchmark.
64
65     Returns:
66     tuple
67         Tuple containing the Sortino ratio of the portfolio and
68         ↪ benchmark.
69     """

```

```

65     _portfolio_downside = round(100 *
        ↳ (np.std(self.portfolio[self.portfolio < 0]) *
        ↳ np.sqrt(252)), 2)
66     _portfolio = round(self.annual_rtn()[0] /
        ↳ _portfolio_downside, 2)[0]
67     _benchmark = None
68     if self.benchmark is not None:
69         _benchmark_downside = round(100 *
            ↳ (np.std(self.benchmark[self.benchmark < 0]) *
            ↳ np.sqrt(252)), 2)
70         _benchmark = round(self.annual_rtn()[1] /
            ↳ _benchmark_downside, 2)[0]
71     return _portfolio, _benchmark
72
73     def annual_VaR(self):
74         """
75         Calculates the annualized Value at Risk (VaR) of the
76         ↳ portfolio and benchmark.
77
78         Returns:
79         tuple
80             Tuple containing the annualized portfolio VaR and
81             ↳ benchmark VaR.
82         """
81         _portfolio = round(self.portfolio.quantile(0.05) *
            ↳ np.sqrt(252), 2)[0]
82         _benchmark = None
83         if self.benchmark is not None:
84             _benchmark = round(self.benchmark.quantile(0.05) *
                ↳ np.sqrt(252), 2)[0]
85     return _portfolio, _benchmark

```

```

86
87 def annual_CVaR(self):
88     """
89     Calculates the annualized Conditional Value at Risk (CVaR)
90     ↪ of the portfolio and benchmark.
91
92     Returns:
93     tuple
94         Tuple containing the annualized portfolio CVaR and
95         ↪ benchmark CVaR.
96     """
97     _portfolio = round(self.portfolio[self.portfolio <
98     ↪ self.portfolio.quantile(0.05)].mean() * np.sqrt(252),
99     ↪ 2)[0]
100     _benchmark = None
101     if self.benchmark is not None:
102         _benchmark = round(self.benchmark[self.benchmark <
103         ↪ self.benchmark.quantile(0.05)].mean() *
104         ↪ np.sqrt(252), 2)[0]
105     return _portfolio, _benchmark
106
107 def worst_drawdown(self):
108     """
109     Calculates the worst drawdown of the portfolio and
110     ↪ benchmark.
111
112     Returns:
113     tuple
114         Tuple containing the worst drawdown of the portfolio and
115         ↪ benchmark.
116     """

```

```

109     _portfolio_cum_rtn = self.portfolio.cumsum()
110     _portfolio_max_rtn = _portfolio_cum_rtn.cummax()
111     _portfolio_drawdown = (_portfolio_cum_rtn -
        ↪ _portfolio_max_rtn) / (1 + _portfolio_max_rtn)
112     _portfolio_worst_drawdown = round(_portfolio_drawdown.min(),
        ↪ 2)[0]
113
114     _benchmark_worst_drawdown = None
115     if self.benchmark is not None:
116         _benchmark_cum_rtn = self.benchmark.cumsum()
117         _benchmark_max_rtn = _benchmark_cum_rtn.cummax()
118         _benchmark_drawdown = (_benchmark_cum_rtn -
        ↪ _benchmark_max_rtn) / (1 + _benchmark_max_rtn)
119         _benchmark_worst_drawdown =
        ↪ round(_benchmark_drawdown.min(), 2)[0]
120
121     return _portfolio_worst_drawdown, _benchmark_worst_drawdown
122
123 def plot_performance(self):
124     """
125     Plots the cumulative returns and drawdown of the portfolio
        ↪ and benchmark.
126     """
127     _portfolio_cum_rtn = self.portfolio.cumsum()
128     _portfolio_max_rtn = _portfolio_cum_rtn.cummax()
129     _portfolio_drawdown = (_portfolio_cum_rtn -
        ↪ _portfolio_max_rtn) / (1 + _portfolio_max_rtn)
130
131     fig, (ax1, ax2) = plt.subplots(2, 1, sharex=True)
132
133     # Plot cumulative returns for portfolio

```



```

134 ax1.plot(_portfolio_cum_rtn.index, _portfolio_cum_rtn,
135         ↪ label='Portfolio Cumulative Returns', color='blue')
136
137 # Plot cumulative returns for benchmark if available
138 if self.benchmark is not None:
139     _benchmark_cum_rtn = self.benchmark.cumsum()
140     ax1.plot(_benchmark_cum_rtn.index, _benchmark_cum_rtn,
141             ↪ label='Benchmark Cumulative Returns', color='green')
142
143 ax1.set_ylabel('Cumulative Returns')
144 ax1.legend()
145
146 # Plot drawdown for portfolio
147 ax2.plot(_portfolio_drawdown.index,
148         ↪ _portfolio_drawdown.values, color='red',
149         ↪ label='Portfolio Drawdown')
150
151 # Plot drawdown for benchmark if available
152 if self.benchmark is not None:
153     _benchmark_max_rtn = _benchmark_cum_rtn.cummax()
154     _benchmark_drawdown = (_benchmark_cum_rtn -
155         ↪ _benchmark_max_rtn) / (1 + _benchmark_max_rtn)
156     ax2.plot(_benchmark_drawdown.index,
157             ↪ _benchmark_drawdown.values, color='orange',
158             ↪ label='Benchmark Drawdown')
159
160 ax2.set_ylabel('Drawdown (%)')
161 ax2.legend()
162
163 # Set x-axis label only on the bottom subplot
164 ax2.set_xlabel('Date')

```

```

158
159     # Customize background
160     fig.set_facecolor('white')
161     ax1.set_facecolor('white')
162     ax2.set_facecolor('white')
163     ax1.grid(color='black', linestyle='--', linewidth=0.2)
164     ax2.grid(color='black', linestyle='--', linewidth=0.2)
165
166     plt.show()
167
168
169     def summary_table(self):
170         """
171         Generates a summary table with various performance metrics
172         → of the portfolio and benchmark.
173
174         Returns:
175         pandas.DataFrame
176         DataFrame containing the summary metrics.
177         """
178         metrics = pd.DataFrame({
179             'portfolio': {'Annualized Returns':
180                 ↪ self.annual_rtn()[0],
181                 'Annualized Standard Deviation':
182                 ↪ self.annual_std()[0],
183                 'Sharpe Ratio': self.sharpe_ratio()[0],
184                 'Sortino Ratio': self.sortino_ratio()[0],
185                 'Annualized VaR 95%':
186                 ↪ self.annual_VaR()[0],
187                 'Annualized CVaR 95%':
188                 ↪ self.annual_CVaR()[0],

```

```

184         'Worst Drawdown':
185             ↪ self.worst_drawdown()[0]},
186     'benchmark': {'Annualized Returns':
187         ↪ self.annual_rtn()[1],
188         'Annualized Standard Deviation':
189             ↪ self.annual_std()[1],
190         'Sharpe Ratio': self.sharpe_ratio()[1],
191         'Sortino Ratio': self.sortino_ratio()[1],
192         'Annualized VaR 95%':
193             ↪ self.annual_VaR()[1],
194         'Annualized CVaR 95%':
195             ↪ self.annual_CVaR()[1],
196         'Worst Drawdown':
197             ↪ self.worst_drawdown()[1]}
198     })
199
200     return metrics

```

## Trading Strategy Execution

```

1  # Define an empty dictionary to store the results
2  strategy = {}
3  portfolio = {}
4
5  for i in range(1, len(inSample)):
6      # Set seed
7      random.seed(2024)
8
9      cointegrated = False
10

```

```

11     while not cointegrated:
12         # Select stocks randomly
13         selected_columns = random.sample(list(inSample[i].columns),
14         ↪ 10)
15         _inSample = inSample[i][selected_columns]
16         _outofSample = outofSample[i][selected_columns]
17
18         # Test for cointegration
19         coint = Johansen(_inSample).n_coint_vectors()
20
21         if coint >= 1:
22             # Run backtest
23             backtest = Backtest(inSample=_inSample,
24             ↪ outofSample=_outofSample, entry_point=1,
25             ↪ stop_loss=3)
26             strategy[i] = pd.DataFrame({"Z-score":
27             ↪ backtest.z_score()}, index=_outofSample.index)
28             portfolio[i] = backtest.port_rtn()
29             cointegrated = True
30         else:
31             # Set backtest to 0 or handle as needed
32             portfolio[i] = pd.DataFrame([0] * len(_outofSample),
33             ↪ index=_outofSample.index)
34
35     # concatenate all DataFrames into a single DataFrame
36     strategy = pd.concat(strategy.values(), axis=0)
37     portfolio = pd.concat(portfolio.values(), axis=0)
38
39     # Plot z-score and signals
40     # Plotting
41     fig, ax = plt.subplots(figsize=(10, 5))

```

```

37 ax.plot(strategy.index, strategy.iloc[:,0], label='Z-score',
    → color='blue', linewidth=1)
38 ax.axhline(y=0, color='black', linestyle='--', linewidth=2,
    → label='ExitPoint')
39 ax.axhline(y=-1, color='green', linestyle='--', linewidth=2,
    → label='Entry Point')
40 ax.axhline(y=1, color='green', linestyle='--', linewidth=2)
41 ax.axhline(y=-3, color='red', linestyle='--', linewidth=2,
    → label='Stop Loss')
42 ax.axhline(y=3, color='red', linestyle='--', linewidth=2)
43 ax.set_ylabel('Values')
44 ax.set_xlabel('Date')
45 ax.set_facecolor('white')
46 ax.grid(color='black', linestyle='--', linewidth=0.3)
47 ax.set_title("Trading Strategy")
48 ax.legend()
49 plt.show()

```

```

1  # Load Ibovespa index
2  ibov = pd.read_excel('Ibovespa.xlsx', skiprows=1, index_col=0,
    → parse_dates=True).dropna()
3  ibov.index.name = 'date'
4
5  # Set as benchmark return
6  benchmark = pd.DataFrame({'Benchmark':
    → np.log(ibov['IBOV']).diff().dropna()})
7  benchmark = benchmark.loc[portfolio.index[0]:portfolio.index[-1]]
8  benchmark.head()

```

```

1  performance = Performance(portfolio=portfolio, benchmark=benchmark)
2  performance.summary_table()

```

```
1 performance.plot_performance()
```

```
1 # Load Fama-French factors
2 ff_factors = pd.read_csv('BR-Fama-French-Factors/FF_Factors.csv',
   ↪ index_col=0, parse_dates=True)
3 ff_factors = ff_factors.loc[portfolio.index[0]:portfolio.index[-1],
   ↪ ['RMRF', 'SMB', 'HML', 'RMW', 'CMA']]
4 ff_factors.head()
```

```
1 # Merge portfolio and ff_factors
2 df = pd.merge(portfolio, ff_factors, left_index=True,
   ↪ right_index=True)
3 df.tail()
```

```
1 # Perform regression
2 X = df.iloc[:,1:]
3 y = df.iloc[:,0]
4 model = sm.OLS(y, X).fit()
5
6 # Display the results
7 model.summary()
```