

Financial Big Data Project:

Implementation of volatility-based strategies

Caterina La Cava (358741), Léo Lagast (290246)

January 29, 2024

Abstract

This project concerns the implementation of different volatility strategies, namely a risk parity strategy, an inverse volatility strategy, and a minimum variance strategy. The data consists of 4 years of tick-by-tick transaction prices on 87 US equity large stocks. Clustering algorithms are used to classify assets into different sub-portfolios, with the objective to apply the trading strategies of above only on these sub-portfolios. The cumulative return of the final portfolio changes as time passes depending on the strategy, the clustering algorithm, and the estimated covariance matrix used in the implementation of the strategies. Louvain clustering and Bahc covariance matrix show the best performance for risk parity and inverse volatility strategies in terms of cumulative return, even outperforming the S&P500.

Introduction

Our objective is to implement known volatility strategies on some obtained sub-portfolio returns: the idea is to extract the data structure from the available data through clustering algorithms to classify assets into groups; weight cluster components with a simple strategy to obtain what we call sub-portfolios; aggregate sub-portfolios in a unique portfolio assigning optimal sub-portfolios weights with known volatility strategies. In the implementation of this idea, there are different decisions to make and each of them could affect the performance of the portfolio's cumulative return over time. Think for example to the choice of the clustering algorithm to use, the way of estimating returns covariance matrices, or to which strategy follow to decide allocation weights in each sub-portfolio. Even the first steps of data pre-processing are essential and sensitive in the pipeline: in fact and most importantly, one of the challenges of the project is to develop our idea above managing a very large quantity of data; thus, we must take decisions on which part of the data available to use, how to clean and treat them and how to solve data specific challenges that we may encounter.

The aim, hence, is to explore as much as we can how different algorithms, assumptions, and decisions could influence the performance of well-known trading strategies when adapting them for day-by-day trading.

We find the topic particularly interesting because it gives us the possibility to approach an allocation problem in different ways, trying combinations of decisions without knowing the "absolute correct" road to follow.

Data set

We obtained a total of 24.5 megabytes of financial data related to US equity transactions regarding 87 different assets. In particular, for each US-listed company available, we work with two kinds of files: bbo and trade files. For both types of files, we have at our disposal 5 years of data, from the beginning of 2004 to the end of 2008, divided into daily files. Daily files correspond to the trading days in each year and, in each of them, are collected intraday data about trade prices, trade volumes (in trade files) and ask prices, bid prices, ask volumes, and bid volumes (in bbo files).

Figure 1 illustrates the organization of the data folder before the pre-processing of the data.

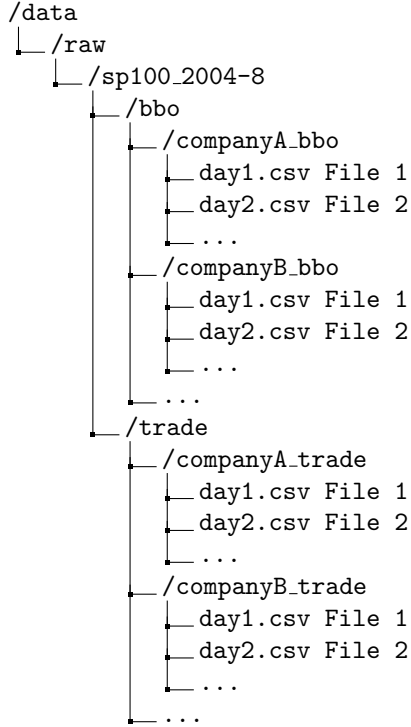


Figure 1: Initial structure of the data folder.

Thus, inside the bbo and trade folders, we find in both 87 folders (one for each asset), that we will call asset folders.

To better work with all the available files, for each asset folder, we proceed as follows.

- Concatenate the daily trade files to obtain a unique data frame whose rows correspond to intraday transactions starting from 2004 and ending in 2008.
- Concatenate the daily bbo files to obtain a unique data frame whose rows correspond to intraday transactions starting from 2004 and ending in 2008.
- Merge the two data frames listed above by datetime.
- Replace the items "()" found in some observations with NaN values.
- Fill NaN values in each column with the previous observation in the same column.
- If present, in each column fill NaN values at the beginning of the data frame with the value of the subsequent observation in the same column.
- Save the final asset data frame in a parquet file and store it in a clean folder within the data folder.

During the cleaning, we noticed that Microsoft (ticker 'MSFT') and Oracle (ticker 'ORCL') have a very high percentage of NaN values; for this reason, we decided to not consider them in our analysis as replacing in any way all these values would have brought to unrealistic prices.

Figure 2 shows how the cleaned data folder looks.

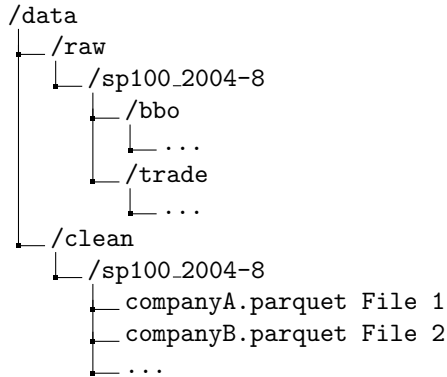


Figure 2: Structure of the data folder with both raw and clean data.

Later, we open the saved parquet files and, for each of them, we replace the observations collected in the same minute by their average to then compute simple returns of the trade price time series. Moreover, we collect the trade price returns of all the assets in one unique dataframe merging them by minute and filling NaN values (deriving from the merging) in each column with the value of the previous observation in the same column. We save the resulting data frame in a parquet file to access it easily (data → clean → returns.parquet).

We explore the returns time series qualitatively by plotting them. As expected, most of the series present high volatility and pronounced peaks starting from the end of 2007 due to the financial crisis; an example is shown in Figure 3.

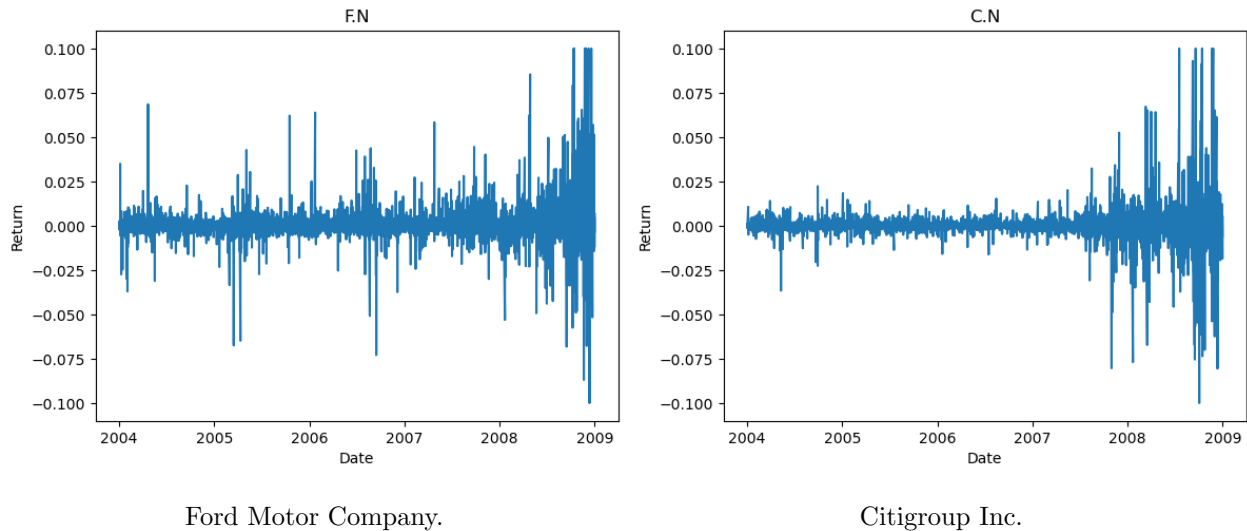
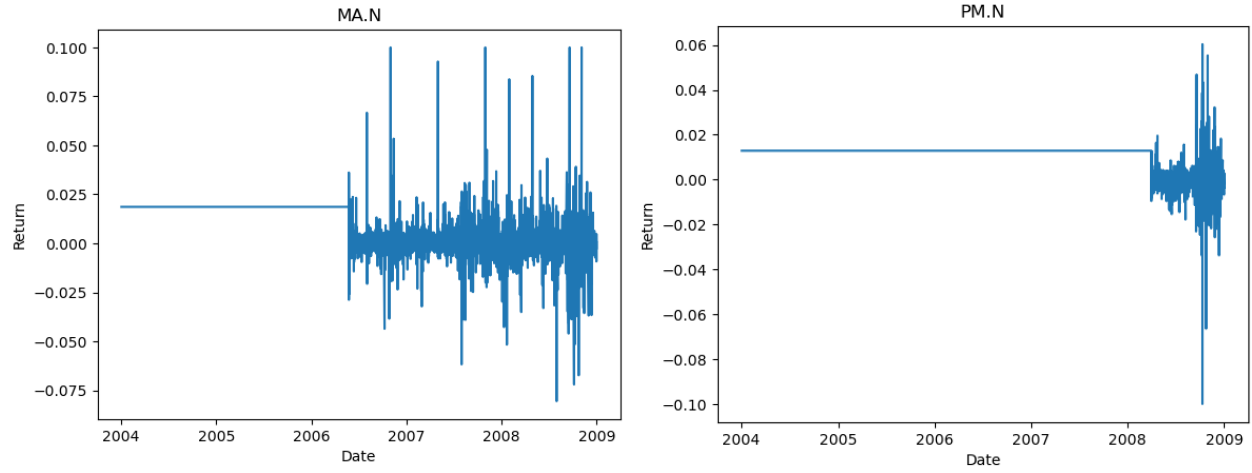


Figure 3: Effects of the financial crisis.

Furthermore, we also notice that some assets have constant returns for a long period, sometimes years; we find in total six assets whose names are available in the Appendix 0.1.

Figure 4 shows an example of the behavior of those time series.



Mastercard Inc.

Philip Morris International.

Figure 4: Bad behavior of some assets.

We then decide to remove bad behavior assets because constant returns cause us problems in computing correct correlation and covariance matrices. We think that the total number of assets, that now amount to 79, is anyway enough for a complete analysis.

Methods

Our idea is to use all the returns collected in the available 5 years time to cluster assets into groups; then construct sub-portfolios whose components are given by assets in the same cluster, build different strategies to weight these sub-portfolios and compare the cumulative returns of the strategies.

We try two different algorithms to cluster assets. First, we use the Louvain clustering algorithm applied to linear asset returns to cluster assets into different groups. We find 6 clusters whose components are shown in the Appendix 0.2. Looking carefully at the companies classified in the same cluster, we notice that the algorithm can divide companies by industry sectors, hence it seems to perform well.

Table 1 shows the main industry sectors represented by the clusters.

Sub-portfolio	Industry
1	Chemicals, manufacturing, technology
2	Healthcare, consumer goods, telecommunications
3	Retail services
4	Financial services
5	Oil producers
6	Energy producers

Table 1: Main industries sectors after performing Louvain clustering.

Then, we use the Marsili Giada clustering algorithm applied to the correlation matrix of asset returns; to estimate the latter, we utilize the clipped correlation matrix which is constructed clipping to a certain threshold the eigenvalues outside of the bulk of the linear correlation matrix. In total, We find 15 clusters which are listed in Appendix ???. We notice immediately the difference in the two clustering algorithms: Marsili Giada groups assets in 15 different clusters; in the first two, there are the majority number of assets, while some of them are composed only of two assets. In the case of Marsili Giada, however, we don't find

a clear pattern in the given clusters as in Louvain. Independently on the cluster algorithm used, in the following, we explain how we decide the allocation within each cluster.

The allocation in each cluster is made dynamically and is given by the minimum variance weights: for a given cluster and at the end of day t , a covariance matrix Σ_t is estimated with the day t returns of the assets in the cluster. New weights are then computed for the sub-portfolio allocation at day $t + 1$. At the end of day $t + 1$, a new covariance matrix Σ_{t+1} is estimated with day $t + 1$ returns and gives the sub-portfolio allocation for day $t + 2$, etc. To estimate Σ_t from returns, we try both the bahc [1] and the clipped covariance matrix; this, then, makes us able to compare results for different combinations of clustering algorithms and covariance matrices.

In the end, we obtain N sub-portfolios (where N is the total number of cluster output of the clustering algorithm) and we develop trading strategies imagining we could trade these sub-portfolios in the market. Moreover, our allocation objective is to decide how to weight these sub-portfolios; the list of the weights associated with the sub-portfolios creates our portfolio.

A more intuitive illustration of how we create our final portfolio is shown in Figure 5.

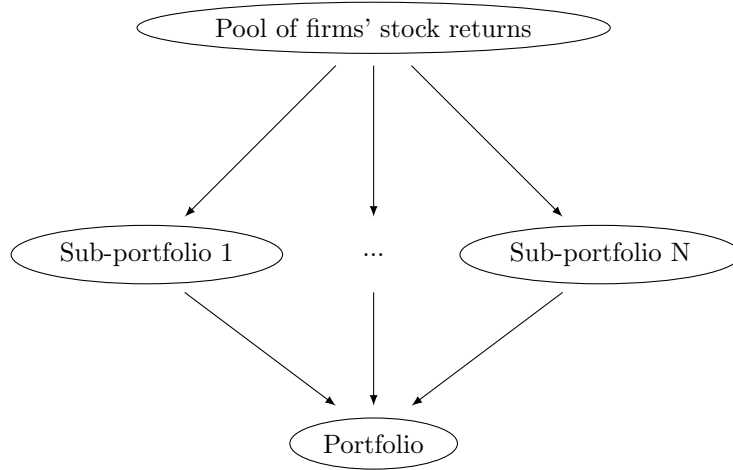


Figure 5: Diagram summarizing the implementation of the portfolio. The components of the N sub-portfolios are obtained by performing a clustering algorithm. The allocation for each sub-portfolio is computed with minimum variance weights. The allocation of every sub-portfolio can vary between risk-parity, inverse volatility, or minimum variance.

Implementation of strategies

Modern portfolio theory, developed by Harry Markowitz in the 1950s, created the foundations for these strategies by introducing the concept of diversification to minimize risk. The starting point for calculating these weights is, in all three cases, an estimate of the covariance matrix of the sub-portfolios, which enables the different strategies determined dynamically to be implemented. The minimum variance portfolio aims to minimize volatility by selecting a fraction of each asset, such that their resulting volatility is minimized while the risk-parity portfolio seeks to balance the contributions of each asset to the portfolio's overall risk, rather than focusing solely on weightings according to their expected return. Risk-parity strategies became popular with Bridgewater's first All-weather funds in the 1990s. This approach is based on the idea that each asset contributes differently to the overall risk of the portfolio, and therefore an equal distribution of risk between assets can lead to better diversification and a reduction in overall portfolio risk. Historically, they seem to handle turbulent phases better than traditional portfolios like 60/40 portfolios. The inverse volatility strategy involves investing more in less volatile assets and less in more volatile assets, to reduce the overall volatility of the portfolio.

Minimum variance weights

The minimum variance weights are computed as follows.

For an n -dimensional time series $\{(r_{1,t}, \dots, r_{n,t})\}_{t=1}^T$ and an estimated covariance matrix Σ_t , the minimum variance weights are given by: the formula

$$w_{mv} = \frac{\Sigma^{-1} \mathbb{1}}{\mathbb{1}^T \Sigma^{-1} \mathbb{1}}.$$

Risk parity weights

It is known that the risk parity strategy tries to give to each asset the allocation that makes it contribute equally to the portfolio aggregate risk. The theory of risk parity can be found in [3] in greater detail, together with examples providing intuition.

Given a portfolio of n assets with weights w , covariance matrix Σ and volatility $\sigma(w) = \sqrt{w^T \Sigma w}$, the Euler Theorem for homogeneous functions implies that $\sigma(w)$ can be decomposed as:

$$\sigma(w) = \sum_{i=1}^n \sigma_i(w),$$

where

$$\sigma_i(w) = \frac{w_i(\Sigma w)_i}{\sqrt{w^T \Sigma w}}.$$

The quantity $\sigma_i(w)$ corresponds to the contribution of asset i to the overall portfolio risk.

Therefore, the risk-parity strategy computes weights such that $\sigma_i(w) = \sigma_j(w)$ for every $i, j \in \{1, \dots, n\}$. The problem hence consists in finding w_{rp} that solves the minimization problem:

$$\arg \min_w \left\{ \sum_{i=1}^n \left(w_i - \frac{\sigma(w)^2}{n(\Sigma w)_i} \right)^2 : \mathbb{1}^T w = 1 \right\}.$$

The implementation of the functions that compute risk-parity weights comes from [2] and has been adapted to our needs.

Inverse volatility weights

The inverse volatility strategy is somehow similar to the risk parity strategy but does not take into account the correlations between assets. If $\Sigma = (\sigma_{ij})_{i,j=1}^n$ is the covariance matrix of n assets, the inverse volatility creates a portfolio with weights w_{iv} , defined as:

$$w_{iv} = \left(\frac{\sigma_{11}}{\text{tr}(\Sigma)}, \dots, \frac{\sigma_{nn}}{\text{tr}(\Sigma)} \right).$$

Results

Figure 6 shows the performance of the 3 considered strategies when applying Louvain clustering and using the Bahc covariance matrix (with $K=1$, $N_{boot}=50$).

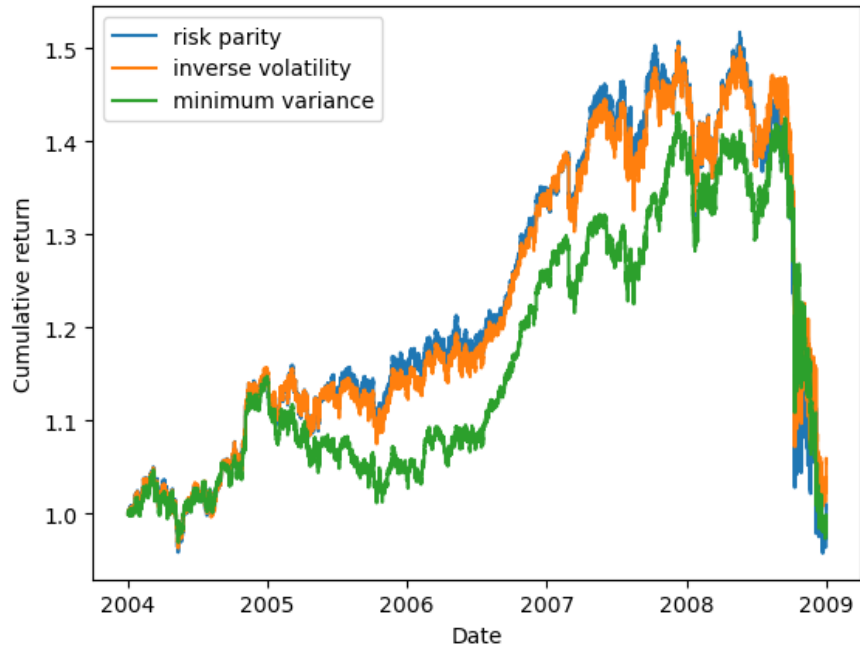


Figure 6: Louvain clustering and bahc covariance matrix.

During the first year all the strategies seem to perform in the same way; then, starting from 2005, risk parity and inverse volatility cumulative returns are similar and both above the minimum variance ones.

Figure 7 shows a comparison between the cumulative returns of the 3 strategies explained above obtained by applying Louvain clustering to classify assets and using the bahc covariance matrix (with $K=1$, $N_{boot}=50$) and the cumulative return of the S&P500.

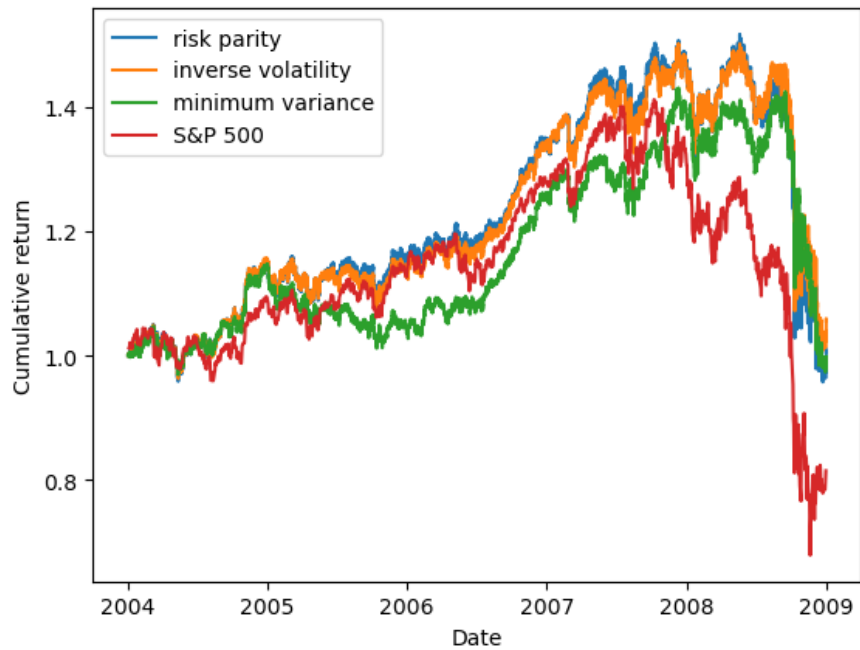


Figure 7: Louvain clustering, Bahc covariance matrix, S&P500.

We notice that the risk parity and inverse volatility strategies outperform the S&P500 almost always, while the minimum variance cumulative return is below the S&P500, except in the last year, during the financial crisis. However, until now, we have not taken into account transaction costs: we are supposed to be able to change our position and trade every day without considering the cost of these trades. This assumption is unrealistic, especially when positions are changed every day. When it comes to considering trading costs, our strategy might not be better than S&P500; it would be interesting to dig into this topic, for example implementing strategies that take into account trading costs.

Figure 8 shows the cumulative returns of the strategies when applying the Marsili Giada clustering algorithm and using the clipped covariance matrix. The S&P500 cumulative return, even if not displayed, is positioned between the minimum variance and the inverse volatility, a bit farther from the inverse volatility than in Figure 7.

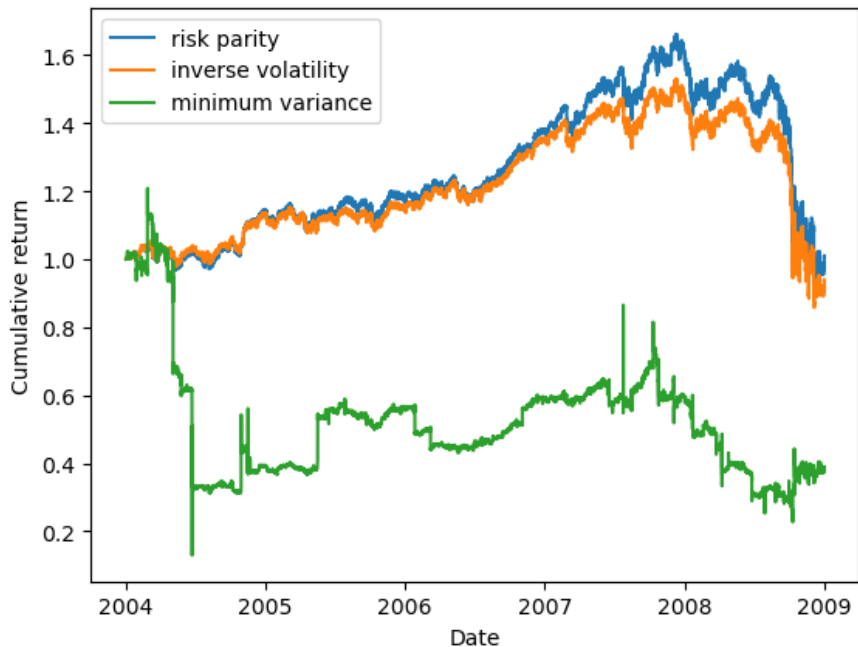


Figure 8: Marsili Giada clustering and clipped covariance matrix.

Here, compared to Figure 6, the minimum variance strategy cumulative return is way below the other strategies; furthermore, it seems less stable than the others, as we can see for example in the middle of 2004 or in the middle of 2007. Moreover, risk parity cumulative return is always above inverse volatility's, reaching also 1.6, something that does not happen in Figure 6.

The comparison between Figure 6 and 8 is an example of how changing the clustering algorithm and how we compute the covariance matrix affect results. We could expect a difference in results especially due to the different output of the two clustering algorithms: Louvain clustering classifies assets into 6 sub-portfolios, while Marsili Giada in 15 sub-portfolios; is quite different.

we try also to obtain the same plot for the combination of Louvain clustering and clipped covariance matrix and, even if the only different input concerning the first case is the covariance matrix, we find a very different plot than Figure 6. We notice that the clipped covariance matrix somehow is less stable than the Bahc one; it gives us different results every time that we run the algorithm and the resulting strategies' cumulative returns fluctuate a lot. For example, in one of the attempts, we find Figure 9.

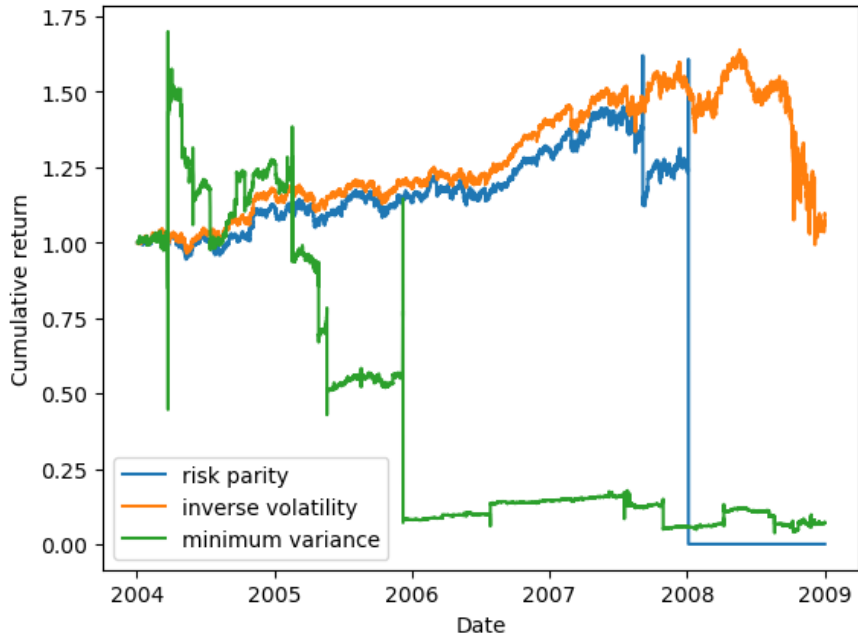


Figure 9: Louvain clustering and clipped covariance matrix.

With the code provided it is possible to run every combination of clustering algorithms and covariance/correlation matrices.

Discussion

The first challenge is to get clean data. Our goal was first to obtain a single data frame containing the linear returns of every asset for every minute in the considered years. The idea behind resampling minute-by-minute is to get rid of the noise and the oscillations of what is inherently contained in tick-by-tick data. This operation lasts several dozens of minutes, even with the use of parallelization and delayed functions.

The second challenge was to decide the destiny of `NaN` values. We decided to apply the radical method of filling a `NaN` with its former value, to replicate the idea that at time t , the best information we have is the last one. As a last resort and especially for the oldest values of the data frame, we applied a backward filling method, i.e. a `NaN` is filled with the first next valid value. There are only a few of them and we estimate that this is not a big issue given the fact that we deal with very small numbers.

These practices may transform drastically the time series of some assets since they can contain large parts of zeros or constants values, but our goal was mostly to have as much possible a data frame with columns of the same length. The main drawback of these fillings appears when we computed the covariance matrix: indeed, if an asset originally had `NaN` values for a whole day, the resulting values of the asset trade price for that day would be either zero or constant and so do its returns. This is problematic as the covariance matrix of returns would then be equal to 0 on the row and the column corresponding to this asset, and therefore would not be invertible. When this phenomenon happened for years since the beginning of the time series, we decided to remove the asset from the analysis (removed assets in 0.1); when it happens for only a small period in the middle of the series (as in the case of Sprint Corporation), we decided to make the asset non-tradable during that period, non considering it as part of the sub-portfolios.

Another problem we encountered is related to the combination of Marsili Giada clustering and Bahc matrices. Above we discussed the case of Marsili Giada clustering and clipped covariance matrix; so, we give the clipped correlation matrix as input of the Marsili Giada algorithm, and in the portfolio construction we

use the clipped covariance matrix. Instead, we tried to put the bahc correlation matrix as input of Marsili Giada and use in the portfolio construction the bahc covariance matrix. In this case, the number of clusters decreased becoming only 3, but then the portfolio construction function did not work, causing different errors. Unfortunately, we had no time to investigate more on the reasons why this happened. However, we were impressed by how different correlation matrices in input could influence the output of Marsili Giada clustering.

There could be a lot of interesting aspects to dig into more in the field of this project's topic, for example, the behavior of the clipped covariance matrix for different types of return time series that maybe could highlight the reasons for the dissimilarities between Figure 6 and Figure 9. We are glad to have the opportunity to work with a very high amount of data and to solve some of the encountered challenges. In conclusion, this project made us able to think about the importance of the decision-making process, to use nonstandard but more suitable methods to treat financial data, and to work with a high volume of information.

Appendix

0.1 List of assets with years of constant returns

Ticker	Company Name
DVN	Devon Energy Corporation
MA	Mastercard Inc.
MS	Morgan Stanley
NOV	Nov Inc.
PM	Philip Morris International Inc.
V	Visa Inc.

Table 2: List of removed assets.

0.2 List of components of each sub-portfolios Louvain clustering

Ticker	Company Name
CAT	Caterpillar Inc.
DD	DuPont de Nemours, Inc.
DOW	Dow Inc.
EMR	Emerson Electric Co.
FDX	FedEx Corporation
GD	General Dynamics Corporation
HON	Honeywell International Inc.
LMT	Lockheed Martin Corporation
MON	Monsanto Company
NKE	NIKE, Inc.
NSC	Norfolk Southern Corporation
RTN	Raytheon Technologies Corporation
UNP	Union Pacific Corporation
UPS	United Parcel Service, Inc.
WY	Weyerhaeuser Company
XRX	Xerox Holdings Corporation

Table 3: Components of sub-portfolio 1.

Ticker	Company Name
ABT	Abbott Laboratories
AVP	Avon Products, Inc.
BAX	Baxter International Inc.
BMY	Bristol-Myers Squibb Company
CL	Colgate-Palmolive Company
HNZ	H.J. Heinz Company
JNJ	Johnson & Johnson
KFT	Kraft Foods Inc.
KO	The Coca-Cola Company
MDT	Medtronic plc
MO	Altria Group, Inc.
MRK	Merck & Co., Inc.
PEP	PepsiCo, Inc.
PFE	Pfizer Inc.
PG	Procter & Gamble Company
S	Sprint Corporation
T	AT&T Inc.
TWX	Time Warner Inc.
VZ	Verizon Communications Inc.

Table 4: Components of sub-portfolio 2.

Ticker	Company Name
BA	Boeing Co.
CVS	CVS Health Corporation
DIS	The Walt Disney Company
HD	The Home Depot, Inc.
HPQ	HP Inc.
IBM	International Business Machines Corporation
LOW	Lowe's Companies, Inc.
MCD	McDonald's Corporation
MMM	3M Company
TGT	Target Corporation
TXN	Texas Instruments Incorporated
UNH	UnitedHealth Group Incorporated
UTX	Raytheon Technologies Corporation
WAG	Walgreen Co.
WMT	Walmart Inc.

Table 5: Components of sub-portfolio 3.

Ticker	Company Name
ALL	Allstate Corporation
AXP	American Express Company
BAC	Bank of America Corporation
BK	The Bank of New York Mellon Corporation
C	Citigroup Inc.
COF	Capital One Financial Corporation
EMC	Emerson Electric Co.
F	Ford Motor Company
GE	General Electric Company
GS	The Goldman Sachs Group, Inc.
JPM	JPMorgan Chase & Co.
MET	MetLife, Inc.
USB	U.S. Bancorp
WFC	Wells Fargo & Co.

Table 6: Components of sub-portfolio 4.

Ticker	Company Name
AA	Alcoa Corporation
APA	Apache Corporation
BHI	Baker Hughes Company
COP	ConocoPhillips
CVX	Chevron Corporation
FCX	Freeport-McMoRan Inc.
HAL	Halliburton Company
OXY	Occidental Petroleum Corporation
SLB	Schlumberger Limited
WMB	The Williams Companies, Inc.
XOM	Exxon Mobil Corporation

Table 7: Components of sub-portfolio 5.

Ticker	Company Name
AEP	American Electric Power Company, Inc.
ETR	Entergy Corporation
EXC	Exelon Corporation
SO	The Southern Company

Table 8: Components of sub-portfolio 6.

0.3 List of components of each sub-portfolios Marsili Giada clustering

Ticker	Company Name
AA	Alcoa Corporation
ABT	Abbott Laboratories
AEP	American Electric Power Company, Inc.
APA	Apache Corporation
BA	Boeing Co.
BAC	Bank of America Corporation
CVS	CVS Health Corporation
EMC	Emerson Electric Co.
EMR	Emerson Electric Co.
ETR	Entergy Corporation
FDX	FedEx Corporation
GE	General Electric Company
HD	The Home Depot Inc.
HPQ	HP Inc.
KFT	Kraft Foods Inc.
LMT	Lockheed Martin Corporation
MDT	Medtronic plc
MMM	3M Company
MO	Altria Group, Inc.
MON	Monsanto Company
MRK	Merck & Co., Inc.
NSC	Norfolk Southern Corporation
SO	Southern Company
VZ	Verizon Communications Inc.
WFC	Wells Fargo & Co.
XOM	Exxon Mobil Corporation

Table 9: Components of sub-portfolio 1.

Ticker	Company Name
ALL	Allstate Corporation
AXP	American Express Company
BHI	Baker Hughes Company
C	Citigroup Inc.
CAT	Caterpillar Inc.
COP	ConocoPhillips
DIS	The Walt Disney Company
F	Ford Motor Company
HAL	Halliburton Company
JPM	JPMorgan Chase & Co.
KO	The Coca-Cola Company
LOW	Lowe's Companies, Inc.
PEP	PepsiCo, Inc.
PFE	Pfizer Inc.
T	AT&T Inc.
TGT	Target Corporation
UNP	Union Pacific Corporation
UPS	United Parcel Service, Inc.
UTX	Raytheon Technologies Corporation
WMT	Walmart Inc.

Table 10: Components of sub-portfolio 2.

Ticker	Company Name
AVP	Avon Products, Inc.
DOW	Dow Inc.
TXN	Texas Instruments Incorporated

Table 11: Components of sub-portfolio 3.

Ticker	Company Name
BAX	Baxter International Inc.
BK	The Bank of New York Mellon Corporation
SLB	Schlumberger Limited
TWX	Time Warner Inc.
WMB	The Williams Companies, Inc.

Table 12: Components of sub-portfolio 4.

Ticker	Company Name
BMJ	Bristol-Myers Squibb Company
JNJ	Johnson & Johnson
S	Sprint Corporation

Table 13: Components of sub-portfolio 5.

Ticker	Company Name
CL	Colgate-Palmolive Company
DD	Du Pont de Nemours, Inc.

Table 14: Components of sub-portfolio 6.

Ticker	Company Name
COF	Capital One Financial Corporation
RTN	Raytheon Technologies Corporation
XXR	Xerox Holdings Corporation

Table 15: Components of sub-portfolio 7.

Ticker	Company Name
CVX	Chevron Corporation
MET	MetLife, Inc.

Table 16: Components of sub-portfolio 8.

Ticker	Company Name
EXC	Exelon Corporation
FCX	Freeport-McMoRan Inc.
GD	General Dynamics Corporation

Table 17: Components of sub-portfolio 9.

Ticker	Company Name
GS	The Goldman Sachs Group, Inc.
IBM	International Business Machines Corporation

Table 18: Components of sub-portfolio 10.

Ticker	Company Name
HNZ	H.J. Heinz Company
OXY	Occidental Petroleum Corporation

Table 19: Components of sub-portfolio 11.

Ticker	Company Name
HON	Honeywell International Inc.
PG	Procter & Gamble Company
USB	U.S. Bancorp

Table 20: Components of sub-portfolio 12.

Ticker	Company Name
MCD	McDonald's Corporation
UNH	UnitedHealth Group Incorporated

Table 21: Components of sub-portfolio 13.

Ticker	Company Name
NKE	NIKE, Inc.

Table 22: Components of sub-portfolio 14.

Ticker	Company Name
WAG	Walgreen Co.
WY	Weyerhaeuser Company

Table 23: Components of sub-portfolio 15.

References

- [1] Bongiorno and Challet (2020). *Package BAHc*. [Online; accessed 28-January-2024]. URL: <https://pypi.org/project/bahc/>.
- [2] fjrodriguez2. *Risk Parity in Python*. [Online; accessed 21-January-2024]. URL: <https://quantdare.com/risk-parity-in-python/>.
- [3] Risk Parity. *Risk Parity — Wikipedia, The Free Encyclopedia*. [Online; accessed 21-January-2024]. URL: https://en.wikipedia.org/wiki/Risk_parity.

Code

```

1 import numpy as np
2 import pandas as pd
3 import os
4 import gzip
5 import tarfile
6 import xlrd
7 import datetime
8 import re
9 import dask
10 import vaex
11 import glob
12 import bahc
13 import community
14 import networkx as nx
15 from scipy.optimize import minimize
16 from numpy import linalg as LA
17 import matplotlib.pyplot as plt
18 import datetime as datetime

```

```

1  # Load and store data
2
3  def get_file_names(folder):
4      """
5          Function that returns an array of strings, where each string is the name of a
6              ↪ file in the folder passed as argument
7      Args:
8          - folder: path
9      Returns:
10         - file_names: array of strings
11     """
12     file_names = []
13     for file in os.listdir(folder):
14         file_names.append(file)
15
16     return file_names
17
18 def create_file(file_name, content):
19     """
20         Creates a file in the data folder, it takes as argument the name of the file
21             ↪ and the content of the file
22     Args:
23         - file_name: name of a file
24         - content: content of the file
25     """
26     with open(file_name, 'w') as file:
27         file.write(content)
28
29 def extract_tar(file_path, output_path):
30     """
31         Extracts the contents of a .tar file to the specified output path.
32     Args:
33         - file_path: name of a file
34         - output_path: path where the content of the file will be stored
35     """
36     try:
37         with tarfile.open(file_path, 'r') as tar:
38             tar.extractall(output_path)
39             print(f"Extraction of {file_path} successful.")
40     except tarfile.TarError as e:
41         print(f"Error extracting {file_path}: {e}")
42
43
44 def extract_csv_gz_file(source_file, destination_directory):
45     """
46         Extracts a .csv.gz file to a specified directory.
47     Args:
48         - source_file: The path to the .csv.gz file to be extracted.

```

```

49     - destination_directory: The directory where the extracted file will be saved.
50     """
51     if not os.path.exists(destination_directory):
52         os.makedirs(destination_directory)
53     try:
54         file_name = os.path.basename(source_file)
55         output_file = os.path.join(destination_directory, os.path.splitext(
56             ↪ file_name)[0])
57         with gzip.open(source_file, 'rb') as f_in, open(output_file, 'wb') as f_out
58             ↪ :
59             f_out.write(f_in.read())
60         print(f"File extracted to: {output_file}")
61     except Exception as e:
62         print(f"Error extracting file: {e}")
63
64 def xl_to_datetime(xltime):
65     """
66     Transforms xltime into an object datetime
67     Args:
68     - xltime: float
69     Returns:
70     - date_time_obj: datetime object
71     """
72     date_value = int(xltime)
73     time_value = (xltime - date_value) * 24 * 60 * 60 # Convert fraction of a day
74                 ↪ to seconds
75     date_tuple = xlrd.xldate_as_tuple(date_value, 0) # 0 for 1900-based date
76                 ↪ system
77     year, month, day, hour, minute, second = date_tuple
78     date_time_obj = datetime.datetime(year, month, day, hour, minute, second) +
79                 ↪ datetime.timedelta(seconds=time_value)
80
81     return date_time_obj
82
83 def convert_to_float(value):
84     """
85     Converts the value to float if it is possible, otherwise it returns nan
86     Args:
87     - value: float
88     Returns:
89     - float_value: or a float or nan
90     """
91     try:
92         float_value = float(value)
93         return float_value if np.isfinite(float_value) else np.nan
94     except (ValueError, TypeError):
95         return np.nan
96

```

```

94
95 def resample_df(df):
96     """
97     Resamples the dataframe df to 1 minute frequency; one apply the function
98     ↪ xl_to_datetime to the column xlttime of merged_df
99
100    Args:
101    - df: dataframe
102    Returns:
103    - df: resampled dataframe
104    """
105    df['datetime'] = df['xlttime'].apply(xl_to_datetime)
106    df['bid-price'] = df['bid-price'].astype(float)
107    df['ask-price'] = df['ask-price'].astype(float)
108    df['bid-volume'] = df['bid-volume'].astype(float)
109    df['ask-volume'] = df['ask-volume'].astype(float)
110
111    #drop the column xlttime
112    df = df.drop(columns=['xlttime'])
113
114    #set the column datetime as index
115    df = df.set_index('datetime')
116    df = df.resample('1T').agg({
117        'bid-price': 'mean',
118        'ask-price': 'mean',
119        'bid-volume': 'sum',
120        'ask-volume': 'sum'
121    })
122
123    return df
124
125 def create_folder(directory_path, folder_name):
126     """
127     Combines directory path and folder name to create the full path for the new
128     ↪ folder
129
130    Args:
131    - directory_path: path
132    - folder_name: name of the folder
133    """
134    new_folder_path = os.path.join(directory_path, folder_name)
135
136    # Create the new folder if it doesn't already exist
137    if not os.path.exists(new_folder_path):
138        os.makedirs(new_folder_path)
139
140 def clean_dataframe(df):
141     """
142     Cleans the dataframe in input by replacing the non-float values with the
143     ↪ previous value

```

```

141     Args:
142     - df: dataframe
143     Returns:
144     - df: cleaned dataframe
145     """
146     for column_name in df.columns[1:]:
147         # Convert the column to numeric, coercing non-numeric values to NaN
148         numeric_column = pd.to_numeric(df[column_name], errors='coerce')
149         mean_value = numeric_column.mean()
150         if isinstance(df[column_name][0], float)==False:
151             df[column_name][0] = mean_value
152         for row in range(1, len(df[column_name])):
153             if isinstance(df[column_name][row], float)==False:
154                 df[column_name][row] = df[column_name][row-1]
155
156     return df
157
158
159 def clean_dataframe_faster(df):
160     """
161     Cleans the dataframe in input by replacing the non-float values with the
162     ↪ previous value
163     Args:
164     - df: dataframe
165     Returns:
166     - df: cleaned dataframe
167     """
168     for column_name in df.columns[1:]:
169         # Convert the column to numeric, coercing non-numeric values to NaN
170         numeric_column = pd.to_numeric(df[column_name], errors='coerce')
171         mean_value = numeric_column.mean()
172
173         # Use vectorized operations to replace non-float values
174         non_float_mask = ~pd.api.types.is_float_dtype(df[column_name])
175
176         df.loc[non_float_mask, column_name] = df[column_name].shift(1)
177         df.loc[0, column_name] = mean_value
178
179     return df
180
181 dask.config.set(scheduler="processes")
182 @dask.delayed
183 def load_trade(filename,
184                tz_exchange="America/New_York",
185                only_non_special_trades=True,
186                only_regular_trading_hours=True,
187                open_time="09:30:00",
188                close_time="16:00:00",
189                merge_sub_trades=True):

```

```

190     """
191     Loads a trade files
192     Args:
193     - filename: name of the file
194     - tz_exchange: timezone
195     - only_non_special_trades: boolean
196     - only_regular_trading_hours: boolean
197     - open_time: string
198     - close_time: string
199     - merge_sub_trades: boolean
200     Returns:
201     - DF: dataframe
202     """
203     try:
204         if re.search('(csv|csv\\.gz)$', filename):
205             DF = pd.read_csv(filename, engine = "pyarrow")
206         if re.search(r'arrow$', filename):
207             DF = pd.read_arrow(filename)
208         if re.search('parquet$', filename):
209             DF = pd.read_parquet(filename)
210     except Exception as e:
211         return None
212     try:
213         DF.shape
214     except Exception as e:
215         print("DF does not exist")
216         print(e)
217         return None
218     if DF.shape[0]==0:
219         return None
220     if only_non_special_trades:
221         DF = DF[DF["trade-stringflag"]=="uncategorized"]
222     DF.drop(columns=["trade-rawflag", "trade-stringflag"], axis=1, inplace=True)
223     DF.index = pd.to_datetime(DF["xltime"], unit="d", origin="1899-12-30", utc=True)
224     DF.index = DF.index.tz_convert(tz_exchange) # .P stands for Arca, which is
        ↪ based at New York
225     DF.drop(columns="xltime", inplace=True)
226     if only_regular_trading_hours:
227         DF=DF.between_time(open_time, close_time) # warning: ever heard e.g.
        ↪ about Thanksgivings?
228     if merge_sub_trades:
229         DF=DF.groupby(DF.index).agg(trade_price=pd.NamedAgg(column='trade-price',
        ↪ , aggfunc='mean'),
230                                     trade_volume=pd.NamedAgg(column='trade-
        ↪ volume', aggfunc='sum'))
231
232     return DF
233
234
235 @dask.delayed

```

```

236 def load_bbo(filename,
237             only_regular_trading_hours=True,
238             merge_sub_trades=True):
239     """
240     Loads a bbo files
241     Args:
242     - filename: name of the file
243     - only_regular_trading_hours: boolean
244     - merge_sub_trades: boolean
245     Returns:
246     - DF: dataframe
247     """
248
249     try:
250         if re.search(r'(csv|csv\.gz)$', filename):
251             DF = pd.read_csv(filename)
252         if re.search(r'arrow$', filename):
253             DF = pd.read_arrow(filename)
254         if re.search(r'parquet$', filename):
255             DF = pd.read_parquet(filename)
256     except Exception as e:
257         return None
258
259     try:
260         DF.shape
261     except Exception as e: # DF does not exist
262         print("DF does not exist")
263         print(e)
264         return None
265
266     if DF.shape[0]==0:
267         return None
268
269     DF.index = pd.to_datetime(DF["xltime"], unit="d", origin="1899-12-30", utc=True)
270     DF.index = DF.index.tz_convert(tz_exchange) # .P stands for Arca, which is
271     ↪ based at New York
272
273     DF.drop(columns="xltime", inplace=True)
274
275     if only_regular_trading_hours:
276         DF=DF.between_time("09:30:00", "16:00:00") # ever heard about
277         ↪ Thanksgivings?
278
279     if merge_sub_trades:
280         DF=DF.groupby(DF.index).last()
281
282     return DF
283
284 @dask.delayed
285 def load_merge_trade_bbo(ticker, date,
286                         dirBase="data/raw/sp100_2004-8/",
287                         suffix="csv.gz",
288                         suffix_save=None,
289                         dirSaveBase="data/clean/sp100_2004-8/events",
290                         saveOnly=False,

```

```

284                 doSave=False
285             ):
286         """
287         Loads and merges the trade and bbo files
288         Args:
289         - ticker: name of the ticker
290         - date: date
291         - dirBase: directory
292         - suffix: suffix
293         - suffix_save: suffix
294         - dirSaveBase: directory
295         - saveOnly: boolean
296         - doSave: boolean
297         Returns:
298         - events: dataframe
299         """
300         file_trade=dirBase+"/"+str(date.date())+"-"+ticker+"-trade
           ↳ ."+suffix
301         file_bbo=file_trade.replace("trade","bbo")
302         trades=load_trade(file_trade)
303         bbos =load_bbo(file_bbo)
304         try:
305             trades.shape + bbos.shape
306         except:
307             return None
308
309         events=trades.join(bbos,how="outer")
310
311         if doSave:
312             dirSave=dirSaveBase+"/"+str(date.date())+"-"+ticker
313             if not os.path.isdir(dirSave):
314                 os.makedirs(dirSave)
315
316             if suffix_save:
317                 suffix=suffix_save
318
319             file_events=dirSave+"/"+str(date.date())+"-"+ticker+"-events"+"."+suffix
320
321             saved=False
322             if suffix=="arrow":
323                 events=vaex.from_pandas(events,copy_index=True)
324                 events.export_arrow(file_events)
325                 saved=True
326             if suffix=="parquet":
327                 # pdb.set_trace()
328                 events.to_parquet(file_events,use_deprecated_int96_timestamps=True)
329                 saved=True
330
331             if not saved:
332                 print("suffix"+suffix+" : format not recognized")

```



```

333
334         if saveOnly:
335             return saved
336
337     return events
338
339
340 def data_to_parquet(ticker):
341     """
342     Loads the data of the ticker in input from the raw folder, cleans it and stores
343     ↪ it in the clean folder
344
345     Args:
346     - ticker: name of the ticker
347     """
348     if ~os.path.exists(f"data/clean/sp100_2004-8/{ticker}.parquet"):
349         trade_files=glob.glob(f"data/raw/sp100_2004-8/trade/{ticker}/*.csv.gz")
350         # we have a name for each trading file that we find in the directory; each
351         ↪ trading file correspond to one
352         # trading day
353         trade_files.sort()
354         allpromises=[load_trade(fn) for fn in trade_files]
355         trades=dask.compute(allpromises)[0]
356         trades=pd.concat(trades)
357
358         bbo_files=glob.glob(f"data/raw/sp100_2004-8/bbo/{ticker}/*.csv.gz")
359         bbo_files.sort()
360         allpromises=[load_bbo(fn) for fn in bbo_files]
361         bbos=dask.compute(allpromises)[0]
362         bbos=pd.concat(bbos)
363
364         events=trades.join(bbos,how="outer")
365
366         # Filling NaNs in 'ask_price' column with the last known value from '
367         ↪ ask_price' column
368         events = events.replace('()', np.nan)
369         events['ask-price'] = events['ask-price'].bfill()
370         events['bid-price'] = events['bid-price'].bfill()
371         events['ask-volume'] = events['ask-volume'].bfill()
372         events['bid-volume'] = events['bid-volume'].bfill()
373         events['ask-price'] = events['ask-price'].ffill()
374         events['bid-price'] = events['bid-price'].ffill()
375         events['ask-volume'] = events['ask-volume'].ffill()
376         events['bid-volume'] = events['bid-volume'].ffill()
377
378         events = events.dropna(subset=['trade_price'])
379         events["bid-price"] = events["bid-price"].values.astype("float")
380         events["bid-volume"] = events["bid-volume"].values.astype("float")
381         events["ask-price"] = events["ask-price"].values.astype("float")
382         events["ask-volume"] = events["ask-volume"].values.astype("float")

```

```

380         events.to_parquet(f"data/clean/sp100_2004-8/{ticker}.parquet")
381
382
383 def process_parquet_files(tickers, trading_returns):
384     """
385     Loads the data of the tickers in input from the clean folder, resamples it and
386     ↪ stores it in the resampled folder
387
388     Args:
389     - tickers: array of tickers
390     - trading_returns: dataframe
391     Returns:
392     - trading_returns: resampled, cleaned dataframe of trading returns
393     """
394     for ticker in tickers:
395         df = pd.read_parquet(f"data/clean/sp100_2004-8/{ticker}")
396         df = df.resample('1T').mean().dropna()
397         df_prices = df.drop(columns=['trade_volume', 'bid-volume', 'ask-volume'])
398         df_returns = (df_prices / df_prices.shift(1) - 1).dropna()
399         trading_returns = pd.concat([trading_returns, df_returns['trade_price'].
400             ↪ to_frame(name= ticker[:-8])], axis=1, join='outer')
401         trading_returns.ffill(inplace=True)
402         trading_returns.bfill(inplace=True)
403
404     for column in trading_returns.columns:
405         trading_returns[column] = np.where(np.abs(trading_returns[column]) > 0.1,
406             ↪ 0.1 * np.sign(trading_returns[column]), trading_returns[column])
407
408     return trading_returns
409
410 # Covariance matrix
411
412 def covariance_matrix(trading_returns_df, if_bahc = False, corr_with_bahc = False):
413     """
414     Computes the covariance matrix of the trading returns dataframe
415     Args:
416     - trading_returns_df: dataframe
417     - bahc: boolean
418     Returns:
419     - cov_matrix: covariance matrix
420     """
421     if if_bahc:
422         if corr_with_bahc:
423             corr = bahc.filterCovariance(trading_returns_df.T.to_numpy(), K = 1,
424                 ↪ Nboot=50, is_correlation=True)
425         else:
426             corr = trading_returns_df.corr()
427         cov_matrix = bahc.filterCovariance(trading_returns_df.T.to_numpy(), K = 1,
428             ↪ Nboot=50)
429
430

```

```

425         return corr, cov_matrix, 0
426     else:
427         # number of timesteps
428         T = trading_returns_df.shape[0]
429         initial_corr = trading_returns_df.corr()
430         # number of assets
431         N = initial_corr.shape[0]
432         if N>T:
433             print("N is bigger than T and risk estimation error may diverge")
434             q = N/T
435             eigenvalues_e, eigenvectors_e = LA.eig(initial_corr)
436             lambda_plus = (1+np.sqrt(q))**2
437             # number of eigenvalues outside of the random bulk
438             outside_bulk_eigenvalues = np.sum(eigenvalues_e>lambda_plus)
439             corr_clipped = eigenvalue_clipping(eigenvalues_e, eigenvectors_e, lambda_plus
440                 → )
441             # from correlation matrix to covariance matrix
442             trading_std = np.array(trading_returns_df.std())
443             cov_matrix = np.outer(trading_std, trading_std) * np.array(corr_clipped)
444
445             return corr_clipped, cov_matrix, outside_bulk_eigenvalues
446
447 def eigenvalue_clipping(lambdas, v, lambda_plus):
448     """
449     Clips the eigenvalues of the correlation matrix to lambda_plus
450     Args:
451     - lambdas: eigenvalues
452     - v: eigenvectors
453     - lambda_plus: threshold
454     Returns:
455     - C_clean: clipped correlation matrix
456     """
457     N=len(lambdas)
458     sel_bulk=lambdas<=lambda_plus
459     N_bulk=np.sum(sel_bulk)
460     sum_lambda_bulk=np.sum(lambdas[sel_bulk])
461     delta=sum_lambda_bulk/N_bulk
462     lambdas_clean=lambdas
463     lambdas_clean[lambdas_clean<=lambda_plus]=delta
464     C_clean=np.zeros((N, N))
465     v_m=np.matrix(v)
466     for i in range(N-1):
467         C_clean=C_clean+lambdas_clean[i] * np.dot(v_m[i,].T, v_m[i,])
468
469     np.fill_diagonal(C_clean,1)
470
471     return C_clean
472
473

```

```

474 # Clustering function to decide which clustering algorithm to use
475
476 def clustering(returns_df, if_marsili_giada = False, if_bahc = False):
477     """
478     Computes the clusters of the returns dataframe
479     Args:
480     - returns_df: dataframe
481     - if_marsili_giada: boolean
482     - if_bahc: boolean
483     Returns:
484     - clusters_constituents: dictionary
485     - sub_portfolios: array of strings
486     - sub_portfolios_returns: dataframe
487     """
488     if if_marsili_giada:
489         corr_matrix, cov_matrix, eigenvalue_clipping = covariance_matrix(returns_df
490             ↪ , if_bahc, True )
491         last_cluster = aggregate_clusters(corr_matrix)
492         s_i = last_cluster['s_i']
493         marsili_clusters = pd.DataFrame()
494         marsili_clusters.index = returns_df.columns
495         marsili_clusters[0] = s_i
496         array_of_numbers = marsili_clusters[0].unique()
497         for i in range(0, len(marsili_clusters[0].unique())):
498             marsili_clusters[0] = np.where(marsili_clusters[0] == array_of_numbers[
499                 ↪ i], i, marsili_clusters[0])
500         clusters = marsili_clusters
501         clusters_dict = marsili_clusters[0].to_dict()
502     else:
503         louvain_clusters = LouvainCorrelationClustering(returns_df)
504         louvain_clusters.index = returns_df.columns
505         clusters = louvain_clusters
506         clusters_dict = louvain_clusters[0].to_dict()
507     sub_portfolios_returns = pd.DataFrame(0.0, index = returns_df.index, columns =
508         ↪ ['sub_pf_{}'.format(i) for i in range(clusters.max().max() + 1)])
509     sub_portfolios = ["sub_pf_{}".format(i+1) for i in range(clusters.max().max() +
510         ↪ 1)]
511     clusters_constituents = {}
512     for letter, value in clusters_dict.items():
513         clusters_constituents.setdefault(value, []).append(letter)
514     return clusters_constituents, sub_portfolios, sub_portfolios_returns
515
516 # Marsili Giada clustering algorithm
517
518 def expand_grid_unique(x, y, include_equals=False):
519     """

```

```

520     Expands the grid of unique values
521     Args:
522     - x: array
523     - y: array
524     - include_equals: boolean
525     Returns:
526     - combinations: array
527     """
528     x = list(set(x))
529     y = list(set(y))
530
531     def g(i):
532         z = [val for val in y if val not in x[:i - include_equals]]
533         if z:
534             return [x[i - 1]] + z
535
536     combinations = [g(i) for i in range(1, len(x) + 1)]
537
538     return [combo for combo in combinations if combo]
539
540
541 def max_likelihood(c, n):
542     """
543     Computes the maximum likelihood
544     Args:
545     - c: float
546     - n: float
547     Returns:
548     - max_likelihood: float
549     """
550     if n > 1:
551         return np.log(n / c) + (n - 1) * np.log((n * n - n) / (n * n - c))
552     else:
553         return 0
554
555 def max_likelihood_list(cs, ns):
556     """
557     Computes the maximum likelihood of a list
558     Args:
559     - cs: dictionary
560     - ns: dictionary
561     Returns:
562     - Lc: dictionary
563     """
564     Lc = {}
565     for x in cs.keys():
566         if ns[x] > 1:
567             Lc[x] = np.log(ns[x] / cs[x]) + (ns[x] - 1) * np.log((ns[x] * ns[x] -
568                 ↪ ns[x]) / (ns[x] * ns[x] - cs[x]))
569         else:

```

```

569         Lc[x] = 0
570
571     return Lc
572
573 def find_max_improving_pair(C, cs, ns, i_s):
574     """
575     Finds the maximum improving pair
576     Args:
577     - C: matrix
578     - cs: dictionary
579     - ns: dictionary
580     - i_s: dictionary
581     Returns:
582     - pair_max_improv: array
583     - Lc_max_impr: float
584     - Lc_old: array
585     """
586     Lc_old = max_likelihood_list(cs, ns)
587     names_cs = list(cs.keys())
588     max_impr = -1e10
589     pair_max_improv = []
590
591     for i in names_cs[:-1]:
592         names_cs_j = names_cs[names_cs.index(i) + 1:]
593         for j in names_cs_j:
594             ns_new = ns[i] + ns[j]
595             i_s_new = i_s[i] + i_s[j]
596             cs_new = np.sum(C[np.ix_(i_s_new, i_s_new)])
597             max_likelihood_new = max_likelihood(cs_new, ns_new)
598             improvement = max_likelihood_new - Lc_old[i] - Lc_old[j]
599
600             if improvement > max_impr:
601                 max_impr = improvement
602                 pair_max_improv = [i, j]
603                 Lc_max_impr = max_likelihood_new
604
605     return {"pair": pair_max_improv, "Lc_new": Lc_max_impr, "Lc_old": [Lc_old[x]
606         ↪ for x in pair_max_improv]}
607
608 def aggregate_clusters(C):
609     """
610     Aggregates the clusters
611     Args:
612     - C: matrix
613     Returns:
614     - last_clusters: dictionary
615     """
616     N = C.shape[0]
617     cs = {i: 1 for i in range(N)}
618     s_i = {i: [i] for i in range(N)}

```

```

618     ns = {i: 1 for i in range(N)}
619     i_s = {i: [i] for i in range(N)}
620     clusters = []
621
622     for i in range(1, N): # hierarchical merging
623         improvement = find_max_improving_pair(C, cs, ns, i_s)
624         Lc_old = improvement['Lc_old']
625         Lc_new = improvement['Lc_new']
626
627         if Lc_new < sum(Lc_old):
628             print("HALF CLUSTER Lc.new > max(Lc.old)")
629
630         if Lc_new <= max(Lc_old):
631             print("Lc.new <= max(Lc.old), exiting")
632             break
633
634         pair = improvement['pair']
635         s_i = [pair[0] if x == pair[1] else x for x in s_i]
636
637         cluster1 = pair[0]
638         cluster2 = pair[1]
639         i_s[cluster1].extend(i_s[cluster2]) # merge the elements of the two
        ➔ clusters
640         del i_s[cluster2] # removes reference to merged cluster2
641
642         ns[cluster1] += ns[cluster2]
643         del ns[cluster2]
644
645         cs[cluster1] = np.sum(C[i_s[cluster1]][:, i_s[cluster1]]) # sums C over
        ➔ the elements of cluster1
646         del cs[cluster2]
647
648         clusters.append({
649             'Lc': max_likelihood_list(cs, ns),
650             'pair_merged': pair,
651             's_i': s_i,
652             'i_s': i_s,
653             'cs': cs,
654             'ns': ns
655         })
656
657     last_clusters = clusters[-1]
658
659     return last_clusters
660
661
662 # Louvain clustering algorithm
663
664 def LouvainCorrelationClustering(R):
665     """

```

```

666     Computes the Louvain clustering of the returns matrix R
667     Args:
668     - R: matrix of returns
669     Returns:
670     - DF: dataframe
671     """
672     N=R.shape[1]
673     T=R.shape[0]
674
675     q=N*1./T
676     lambda_plus=(1.+np.sqrt(q))*2
677
678     C=R.corr()
679     lambdas, v = LA.eigh(C)
680     C_s=compute_C_minus_C0(lambdas,v,lambda_plus)
681
682     mygraph= nx.from_numpy_matrix(np.abs(C_s))
683     partition = community.community_louvain.best_partition(mygraph)
684
685     DF=pd.DataFrame.from_dict(partition,orient="index")
686
687     return(DF)
688
689
690 def compute_C_minus_C0(lambdas,v,lambda_plus,removeMarketMode=True):
691     """
692     Computes the correlation matrix C minus C0
693     Args:
694     - lambdas: eigenvalues
695     - v: eigenvectors
696     - lambda_plus: threshold
697     - removeMarketMode: boolean
698     Returns:
699     - C_clean: correlation matrix
700     """
701     N=len(lambdas)
702     C_clean=np.zeros((N, N))
703
704     order = np.argsort(lambdas)
705     lambdas,v = lambdas[order],v[:,order]
706
707     v_m=np.matrix(v)
708
709     for i in range(1*removeMarketMode,N):
710         if lambdas[i]>lambda_plus:
711             C_clean=C_clean+lambdas[i] * np.dot(v_m[:,i],v_m[:,i].T)
712
713     return C_clean
714
715

```



```

716 # Portfolios construction
717 # Part of the following code comes from https://quantdare.com/risk-parity-in-python
    ↪ / and has
718 # been adapted to our needs.
719
720 def portfolios_construction(clusters_constituents, returns, dates,
    ↪ sub_portfolios_returns, if_bahc = False):
721     """
722     Constructs the portfolios
723     Args:
724     - clusters_constituents: dictionary
725     - returns: dataframe
726     - dates: array
727     - sub_portfolios_returns: dataframe
728     - if_bahc: boolean
729     Returns:
730     - portfolio_cumprod: cumulative returns dataframe
731     - portfolio_return: portfolio returns dataframe
732     """
733     for cluster in clusters_constituents.keys():
734         df_without_0, stocks_to_be_removed = remove_bad_columns(returns[
    ↪ clusters_constituents[cluster]].loc[dates[0]])
735         corr_clipped, cov_matrix, outside_bulk_eigenvalues = covariance_matrix(
    ↪ df_without_0, if_bahc)
736         w_mv = compute_mv_weights(cov_matrix)
737
738         for date in dates[1:]:
739             clusters_constituents_without_stocks_to_be_removed = [name for name in
    ↪ clusters_constituents[cluster] if name not in
    ↪ stocks_to_be_removed]
740             df_without_0, stocks_to_be_removed = remove_bad_columns(returns[
    ↪ clusters_constituents[cluster]].loc[date])
741             corr_clipped, cov_matrix, outside_bulk_eigenvalues = covariance_matrix(
    ↪ df_without_0, if_bahc)
742
743             series = pd.Series(np.dot(w_mv.T, returns[
    ↪ clusters_constituents_without_stocks_to_be_removed].loc[date].T)
    ↪ )
744             series.index = sub_portfolios_returns.loc[date].index
745             sub_portfolios_returns[f'sub_pf_{cluster}'].loc[date] = series
746
747             w_mv = compute_mv_weights(cov_matrix)
748
749     # sub_portfolios_returns = pd.DataFrame(np.where(np.abs(sub_portfolios_returns)
    ↪ > 0.5, 0.5 * np.sign(sub_portfolios_returns), sub_portfolios_returns))
750     date = dates[1]
751     portfolio_return = pd.DataFrame(0.0, index = returns.index, columns = ['rp_pf',
    ↪ 'iv_pf', 'mv_pf'])
752     rp_weights = pd.DataFrame(0.0, index = dates, columns = ["sub_pf_{}".format(i)
    ↪ for i in range(max(list(clusters_constituents.keys())) + 1)])

```

```

753
754     for date in dates[1:]:
755         df_without_0, stocks_to_be_removed = remove_bad_columns(
756             ↪ sub_portfolios_returns.loc[date])
757         corr_clipped, cov_matrix, outside_bulk_eigenvalues = covariance_matrix(
758             ↪ sub_portfolios_returns.loc[date], if_bahc)
759         w_rp = compute_rp_weights(cov_matrix)
760         w_iv = compute_iv_weights(cov_matrix)
761         w_mv = compute_mv_weights(cov_matrix)
762
763         series_rp = pd.Series(np.dot(w_rp.T, sub_portfolios_returns.loc[date].T))
764         series_rp.index = portfolio_return.loc[date].index
765         portfolio_return['rp_pf'].loc[date] = series_rp
766
767         series_iv = pd.Series(np.dot(w_iv.T, sub_portfolios_returns.loc[date].T))
768         series_iv.index = portfolio_return.loc[date].index
769         portfolio_return['iv_pf'].loc[date] = series_iv
770
771         series_mv = pd.Series(np.dot(w_mv.T, sub_portfolios_returns.loc[date].T))
772         series_mv.index = portfolio_return.loc[date].index
773         portfolio_return['mv_pf'].loc[date] = series_mv
774
775         series_w_rp = pd.Series(w_rp)
776         series_w_rp.index = rp_weights.loc[date].index
777         rp_weights.loc[date] = series_w_rp
778
779     for column in portfolio_return.columns:
780         portfolio_return[column] = np.where(np.abs(portfolio_return[column]) > 0.3,
781             ↪ 0.3 * np.sign(portfolio_return[column]), portfolio_return[column])
782
783     portfolio_cumprod = (portfolio_return + 1).cumprod()
784
785     for column in portfolio_cumprod.columns:
786         portfolio_cumprod[column] = np.where(np.abs(portfolio_cumprod[column]) >
787             ↪ 1.7, 1.7 * np.sign(portfolio_cumprod[column]), portfolio_cumprod[
788             ↪ column])
789
790     return portfolio_cumprod, portfolio_return
791
792 def compute_mv_weights(cov):
793     """
794     Computes the weights of a minimal variance portfolio given the covariance
795     ↪ matrix
796
797     Args:
798     - cov: covariance matrix
799
800     Returns:
801     - w: weights
802     """

```

```

797     eigenvalues, eigenvectors = np.linalg.eig(cov)
798     D = np.diag(eigenvalues)
799     inv_D = np.linalg.inv(D)
800     inv_cov = np.dot(np.dot(eigenvectors, inv_D), eigenvectors.T)
801     w = inv_cov.dot(np.ones(len(cov))) / np.ones(len(cov)).T.dot(inv_cov).dot(np.
        ↪ ones(len(cov)))
802
803     return w
804
805
806 def compute_iv_weights(cov):
807     """
808     Computes the weights of an inverse volatility portfolio given the covariance
        ↪ matrix
809     Args:
810     - cov: covariance matrix
811     Returns:
812     - w: weights
813     """
814     w = 1 / np.diag(cov)
815     w = w / w.sum()
816
817     return w
818
819
820 def compute_vol(w, cov):
821     """
822     Computes the volatility of a portfolio given the weights and the covariance
        ↪ matrix
823     Args:
824     - w: weights
825     - cov: covariance matrix
826     Returns:
827     - pf_risk: portfolio risk
828     """
829     pf_risk = np.dot(w, cov).dot(w.T)
830
831     return pf_risk
832
833 def compute_assets_risk_contribution(w, cov):
834     """
835     Computes the contribution of each asset to the risk of the whole portfolio
836     Args:
837     - w: weights
838     - cov: covariance matrix
839     Returns:
840     - assets_risk_contribution: array of assets risk contribution
841     """
842     pf_risk = compute_vol(w, cov)
843     assets_risk_contribution = np.multiply(w.T, np.dot(cov, w.T)) / pf_risk

```

```

844
845     return assets_risk_contribution
846
847
848 def objective_function(w, args):
849     """
850     Objective function to minimize
851     Args:
852     - w: weights
853     - args: arguments
854     Returns:
855     - error: error
856     """
857     cov = args[0]
858     assets_risk_budget = args[1]
859     w = np.matrix(w)
860     pf_risk = compute_vol(w, cov)
861     assets_risk_contribution = compute_assets_risk_contribution(w, cov)
862     assets_risk_target = np.asmatrix(np.multiply(pf_risk, assets_risk_budget))
863     error = sum(np.square(assets_risk_contribution - assets_risk_target.T))[0, 0]
864
865     return error
866
867 def compute_rp_weights(cov, constraint_type = 'long'):
868     """
869     Computes the weights of a risk parity portfolio given the covariance matrix
870     Args:
871     - cov: covariance matrix
872     - constraint_type: type of constraint
873     Returns:
874     - w: weights
875     """
876     initial_weights = np.ones(cov.shape[0]) * (1.0 / cov.shape[0])
877     assets_risk_budget = np.ones(cov.shape[0]) * (1.0 / cov.shape[0])
878     if constraint_type == 'long':
879         #constraints for long-only portfolio. sum of weights = 1
880         constraints = ({'type': 'eq', 'fun': lambda x: np.sum(x) - 1.0}, {'type': '
            ↳ ineq', 'fun': lambda x: x})
881     elif constraint_type == 'long_short':
882         #constraints for long-short portfolio. sum of weights = 0
883         constraints = ({'type': 'eq', 'fun': lambda x: np.sum(x)})
884     else:
885         raise ValueError('Unknown constraint_type. Only possible values for the
            ↳ get_risk_parity_weights constraint_type argument are "long" and "
            ↳ long_short"')
886
887     sol = minimize(fun=objective_function, x0=initial_weights, args=[cov,
            ↳ assets_risk_budget], method='SLSQP', constraints=constraints, tol=1e-15,
            ↳ options={'disp': False})
888     w = sol.x

```

```

889
890     return w
891
892
893 def remove_bad_columns(dataframe):
894     """
895     Removes the columns of a dataframe that are constant, NaN or infinite
896     Args:
897     - dataframe: dataframe
898     Returns:
899     - new_df: new dataframe
900     - list(zero_columns): list of removed columns
901     """
902     for column in dataframe.columns:
903         if dataframe[column].nunique() == 1:
904             dataframe[column] = 0
905     dataframe.fillna(0, inplace=True)
906     dataframe.replace([np.inf, -np.inf], 0, inplace=True)
907     zero_columns = dataframe.columns[dataframe.eq(0).all()]
908     new_df = dataframe.loc[:, ~dataframe.columns.isin(zero_columns)]
909
910     return new_df, list(zero_columns)
911
912
913 def SP500_returns(path):
914     """
915     Computes the returns of the S&P500 index
916     Args:
917     - path: path
918     Returns:
919     - portfolio_cumprod_SP500: cumulative returns dataframe
920     """
921     df_sp500 = pd.read_csv(path)
922     df_reversed = df_sp500.iloc[::-1]
923     df_reversed.reset_index(inplace=True)
924     df_reversed.drop(columns=['index'], inplace=True)
925     df_reversed['Date'] = pd.to_datetime(df_reversed['Date'])
926     df_reversed.set_index('Date', inplace=True)
927     df = df_reversed['Price'].to_frame()
928     df['Price'] = df['Price'].str.replace(',', '').astype(float)
929     df = (df / df.shift(1) - 1).dropna()
930     df.rename(columns={'Price': 'Returns_S&P500'}, inplace=True)
931     portfolio_cumprod_SP500 = (df + 1).cumprod()
932     # save dataframe into a parquet file
933     portfolio_cumprod_SP500.to_parquet('data/clean/SP500returns.parquet', engine='
        ↪ pyarrow')
934
935     return portfolio_cumprod_SP500

```

```

1 # stock_tickers contains the name of all the files/stocks

```

```

2 stock_tickers_bbo = get_file_names('data/raw/sp100_2004-8/bbo')
3 stock_tickers_bbo.remove('.DS_Store')
4 for file_name_bbo in stock_tickers_bbo:
5     file_path_bbo = f"data/raw/sp100_2004-8/bbo/{file_name_bbo}/{file_name_bbo}_bbo
        ↪ .tar"
6     output_path_bbo = f"data/raw/sp100_2004-8/bbo/{file_name_bbo}/"
7     extract_tar(file_path_bbo, output_path_bbo)

1 stock_tickers_trade = get_file_names('data/raw/sp100_2004-8/trade')
2 for file_name_trade in stock_tickers_trade:
3     file_path_trade = f"data/raw/sp100_2004-8/trade/{file_name_trade}/{
        ↪ file_name_trade}_trade.tar"
4     output_path_trade = f"data/raw/sp100_2004-8/trade/{file_name_trade}/"
5     extract_tar(file_path_trade, output_path_trade)

1 tickers = get_file_names("data/raw/sp100_2004-8/trade/")
2 tickers.remove('MSFT.O')
3 tickers.remove('ORCL.N')
4 for ticker in tickers:
5     data_to_parquet(ticker)

1 trading_returns = pd.DataFrame()
2 tickers = get_file_names("data/clean/sp100_2004-8/")
3 tickers.remove('DVN.N.parquet')
4 tickers.remove('MA.N.parquet')
5 tickers.remove('MS.N.parquet')
6 tickers.remove('NOV.N.parquet')
7 tickers.remove('PM.N.parquet')
8 tickers.remove('V.N.parquet')
9 trading_returns_df = process_parquet_files(tickers, trading_returns)
10 trading_returns_df.to_parquet('data/clean/returns.parquet', engine='pyarrow')

1 # run this cell only if you want to recreate the file already stored with the name
    ↪ "portfolio_cumprod_louvain_bahc.parquet"
2 returns = pd.read_parquet(f"data/clean/returns.parquet")
3 dates = returns.index.strftime("%Y-%m-%d").unique()
4 clusters_constituents, sub_portfolios, sub_portfolios_returns = clustering(returns,
    ↪ if_marsili_giada=False)
5 portfolio_cumprod = portfolios_construction(clusters_constituents, returns, dates,
    ↪ sub_portfolios_returns, if_bahc=True)
6 portfolio_cumprod.to_parquet('data/clean/portfolio_cumprod_louvain_bahc.parquet')

1 portfolio_cumprod = pd.read_parquet('data/clean/portfolio_cumprod_louvain_bahc.
    ↪ parquet')
2 plt.plot(portfolio_cumprod['rp_pf'], label = 'risk_parity')
3 plt.plot(portfolio_cumprod['iv_pf'], label = 'inverse_volatility')
4 plt.plot(portfolio_cumprod['mv_pf'], label = 'minimum_variance')
5 plt.xlabel('Date')
6 plt.ylabel('Cumulative_return')
7 plt.legend()

```

```

8 plt.show()

```

```

1 portfolio_cumprod = pd.read_parquet('data/clean/portfolio_cumprod_louvain_bahc.
    ↪ parquet')
2 path = "data/raw/sp100_2004-8/S&P_500_Historical_Data.csv"
3 portfolio_cumprod_SP500 = SP500_returns(path)
4 plt.plot(portfolio_cumprod['rp_pf'], label = 'risk_parity')
5 plt.plot(portfolio_cumprod['iv_pf'], label = 'inverse_volatility')
6 plt.plot(portfolio_cumprod['mv_pf'], label = 'minimum_variance')
7 plt.plot(portfolio_cumprod_SP500, label = 'S&P_500')
8 plt.xlabel('Date')
9 plt.ylabel('Cumulative_return')
10 plt.legend()
11 plt.show()

```

```

1 # run this cell only if you want to recreate the file already stored with the name
    ↪ "portfolio_cumprod_marsili_clipped.parquet"
2 returns = pd.read_parquet(f"data/clean/returns.parquet")
3 dates = returns.index.strftime("%Y-%m-%d").unique()
4 clusters_constituents, sub_portfolios, sub_portfolios_returns = clustering(returns,
    ↪ if_marsili_giada=True)
5 portfolio_cumprod, portfolio_return = portfolios_construction(clusters_constituents
    ↪ , returns, dates, sub_portfolios_returns, if_bahc=False)
6 portfolio_cumprod.to_parquet('data/clean/portfolio_cumprod_marsili_clipped.parquet'
    ↪ )

```

```

1 portfolio_cumprod = pd.read_parquet('data/clean/portfolio_cumprod_marsili_clipped.
    ↪ parquet')
2 plt.plot(portfolio_cumprod['rp_pf'], label = 'risk_parity')
3 plt.plot(portfolio_cumprod['iv_pf'], label = 'inverse_volatility')
4 plt.plot(portfolio_cumprod['mv_pf'], label = 'minimum_variance')
5 plt.xlabel('Date')
6 plt.ylabel('Cumulative_return')
7 plt.legend()
8 plt.show()

```

```

1 portfolio_cumprod = pd.read_parquet('data/clean/portfolio_cumprod_marsili_clipped.
    ↪ parquet')
2 path = "data/raw/sp100_2004-8/S&P_500_Historical_Data.csv"
3 portfolio_cumprod_SP500 = SP500_returns(path)
4 plt.plot(portfolio_cumprod['rp_pf'], label = 'risk_parity')
5 plt.plot(portfolio_cumprod['iv_pf'], label = 'inverse_volatility')
6 plt.plot(portfolio_cumprod['mv_pf'], label = 'minimum_variance')
7 plt.plot(portfolio_cumprod_SP500, label = 'S&P_500')
8 plt.xlabel('Date')
9 plt.ylabel('Cumulative_return')
10 plt.legend()
11 plt.show()

```

```

1  # run this cell only if you want to recreate the file already stored with the name
    ↪ "portfolio_cumprod_louvain_clipped.parquet"
2  returns = pd.read_parquet(f"data/clean/returns.parquet")
3  dates = returns.index.strftime("%Y-%m-%d").unique()
4  clusters_constituents, sub_portfolios, sub_portfolios_returns = clustering(returns,
    ↪ if_marsili_giada=False)
5  portfolio_cumprod, portfolio_return = portfolios_construction(clusters_constituents
    ↪ , returns, dates, sub_portfolios_returns, if_bahc=False)
6  portfolio_cumprod.to_parquet('data/clean/portfolio_cumprod_louvain_clipped.parquet'
    ↪ )

```

```

1  portfolio_cumprod = pd.read_parquet('data/clean/portfolio_cumprod_louvain_clipped.
    ↪ parquet')
2  plt.plot(portfolio_cumprod['rp_pf'], label = 'risk_parity')
3  plt.plot(portfolio_cumprod['iv_pf'], label = 'inverse_volatility')
4  plt.plot(portfolio_cumprod['mv_pf'], label = 'minimum_variance')
5  plt.xlabel('Date')
6  plt.ylabel('Cumulative_return')
7  plt.legend()
8  plt.show()

```

```

1  portfolio_cumprod = pd.read_parquet('data/clean/portfolio_cumprod_louvain_clipped.
    ↪ parquet')
2  path = "data/raw/sp100_2004-8/S&P_500_Historical_Data.csv"
3  portfolio_cumprod_SP500 = SP500_returns(path)
4  plt.plot(portfolio_cumprod['rp_pf'], label = 'risk_parity')
5  plt.plot(portfolio_cumprod['iv_pf'], label = 'inverse_volatility')
6  plt.plot(portfolio_cumprod['mv_pf'], label = 'minimum_variance')
7  plt.plot(portfolio_cumprod_SP500, label = 'S&P_500')
8  plt.xlabel('Date')
9  plt.ylabel('Cumulative_return')
10 plt.legend()
11 plt.show()

```
