

## ✓ PORTFOLIO MANAGEMENT GWP3

---

### INTRODUCTION:

In this project, we are going to use the portfolios used in our last project GWP2. For better portfolio management, we are going to find the portfolio with the best risk-reward relationship and do improvements on the portfolio using methods like denoising, clustering and backtesting. Multiple improvements are also used on a single portfolio. With these different improvements, we are going to find the performance of our portfolio. These improvements could give advantage to our portfolio. We are going to analyze the performance of our portfolio with these different improvements in this portfolio optimization process. The project belows intends to bring improvement for asset allocation for best portfolio in the previous allocation under project 2 using the denoising, clustering and backtesting.

### Step 1: Theory about denoising, clustering and backtesting.

#### IMPROVEMENTS USING DENOISING:

Denoising techniques are important in improving the accuracy and reliability of the models by removing the meaningful signals from noise. Mostly the estimation errors are mentioned as noise. There are various denoising methods. De Prado also introduced methods to improve the quality of data using denoising.

#### FEATURES:

- Denoising techniques are very useful and robust against errors and outliers.
- The Principal Component Analysis (PCA) is a denoising method that helps in decomposing the covariance matrix and also isolates the principal components. The principal components are actually the meaningful signals. By focusing on the principal components, PCA reduces the dimensionality and filters out the noise. This improves the data quality and thus lowers the amount of data without losing much of the information.
- Denoising helps in replacing the eigenvalues of the eigenvectors classified as random by Marcenko-Pastur with a constant eigenvalue[1]. This helps in reducing the noise present in the correlation matrix. The important signals are kept while the noise is removed.
- Random Matrix Theory (RMT) with Marcenko-Pastur distribution (MPD), helps in finding the difference between signal and noise in the eigenvalues of the covariance matrix. This helps in finding the non-random eigenvalues.
- Constant Residual Eigenvalue Method (CREM) reduces noise by replacing the noisy eigenvalues with their mean. This helps in getting a more stable covariance matrix for

portfolio optimization.

- By introducing a shrinkage coefficient, Targeted shrinkage can be used for the adjustment of the denoising degree. This method allows control of noise reduction and in keeping the important signals. It reduces noise and gives optimal model performance. Targeted shrinkage and CREM are introduced by De Prado.
- MIC-EMD is an adaptive denoising method that can help in removing the noise present in input features. This is done based on the nonlinear relationship between the target variable and the input features[2].

### **BENEFITS:**

- The denoising process reduces the noise in the covariance matrix. This helps in getting more accurate portfolio optimization. This improves the risk-return ratio of the portfolios.
- This helps in finding and keeping the principal components that represent the important signals. Thus helping to make better investment decisions.
- Denoising helps in preventing overfitting by removing the noisy components. The noise reduction process decreases the estimation error in the covariance matrix.
- With the targeted shrinkage method, the degree of denoising can be adjusted. In a variance-covariance matrix after denoising, we can find the uncorrelated assets easily. This helps in getting a more effective diversified portfolio and in reducing the overall portfolio risk.
- Denoising helps in reducing risk and increasing the returns. Portfolio risk can be found more accurately. The accuracy of the model can also be increased. Denoising helps in reducing the importance of non-stationary.
- It also helps in portfolio optimization by finding the optimal asset allocation.
- Stable portfolio performance can be achieved due to denoising.

### **IMPROVEMENTS USING CLUSTERING:**

Based on the similarity, assets are organized into groups. This is called clustering. To make portfolios more diversified and manage risk, clustering can be used.

### **FEATURES:**

- Assets can be grouped based on similarity. These similar characteristics include volatility, performance or market cap and several others. Clustering is an unsupervised learning method.
- To measure the similarity between assets, different distance metrics can be used. Some of the common distance metrics are correlation distances and Euclidean distances. A distance metric must satisfy the three conditions: non-negativity, symmetry, and the triangle inequality.

- Hierarchical Clustering arranges assets into a hierarchy. Hierarchical Risk Parity (HRP) is risk parity with hierarchical clustering. This can improve portfolio optimization.
- Gaussian Mixture Model is also a clustering method that forms groups by using mean and variance.
- Clustering can handle large datasets easily.
- K-means algorithm can also be used for clustering.

### **BENEFITS:**

- Grouping the assets with similar risk together, clustering can help in managing the risk.
- Using Hierarchical Clustering, the clustered portfolio gives higher risk-adjusted returns, better standard deviation rates, and improved Sharpe-ratios.
- This method also helps in increasing the robustness of the model.
- Because of the grouping of similar assets, diversified portfolios can be made. Diversification reduces the importance of individual asset volatility on the overall portfolio.
- Mutual Information helps in finding the reduction of the uncertainty. Clustering helps in finding both the linear and non-linear relationships between assets.
- Hierarchical Risk Parity helps in reducing the noise in the data. This helps in getting a more stable portfolio performance.
- For finding the best optimal asset allocation, clustering can be done by grouping assets based on different characteristics. This helps in risk management.

### **IMPROVEMENTS USING BACKTESTING:**

Backtesting is a method, which helps in finding the performance of models. This is used in portfolio management. We can do this backtesting by simulating how a strategy would have performed in the past. This helps in finding the gains and losses of the strategy.

### **FEATURES:**

- Backtesting can be done by finding the accuracy of the predictions we got in the model. Mean Squared Error (MSE) is one of the metrics.
- These metrics help in finding how much the model's predictions are different from the actual historical data.
- Sharpe Ratio is also one of the backtesting methods. This has both return and risk. The performance of the strategy used in backtesting depending on risk can be calculated with this.
- Value at risk is also one of the backtesting methods.
- Walk-forward backtesting helps in testing the model on out-of-sample periods. This method helps in finding the robustness of the model. This can also help in finding if the

strategy can be used during different market conditions. Cross-validation backtesting uses k-fold cross-validation. This method also helps in finding the robustness of the model.

### **BENEFITS:**

- Backtesting helps in finding if the model is overfitting or underfitting. On unseen data, the model is tested to find the performance.
- In backtesting, the historical data is used to check the strategy performance. The robustness of the model can be found during backtesting.
- During backtesting, the risk of the model can be found using the risk metrics. This can help in risk management.
- The potential risks can be found during the backtesting and then we can reduce the risk of the strategy.
- The model's hyperparameters can be fine-tuned during backtesting. This can help in improving the performance of the strategy.
- The backtesting helps in finding if the strategy can be used or not.
- During backtesting, by checking on the performance of strategy, we can find where the model performs better and the strategy's weakness.
- During different market conditions, if the strategy can be used or not is found during the backtesting.
- During backtesting, we can find the parameters that need to be optimized to improve the performance of strategy.
- The risk-reward ratio can be found during this and thus we can help the strategy have a better risk-reward ratio.

## ✓ Step 2: The improvement on the best portfolio from project 2

### Yahoo Finance

```
! pip install -q yfinance
```

### Packages

```
import pandas as pd
import yfinance as yf
import numpy as np
```

### Top 100 companies

```

import yfinance as yf
import pandas as pd

# Fetch the list of S&P 500 tickers
sp500_tickers = pd.read_html('https://en.wikipedia.org/wiki/List_of_S%26P_500_companies')[0]['Symbol'].t

# Fetch market capitalization data for the S&P 500 companies
market_caps = {}
for ticker in sp500_tickers:
    try:
        company = yf.Ticker(ticker)
        market_cap = company.info['marketCap']
        market_caps[ticker] = market_cap
    except KeyError:
        print(f"Market cap data not found for {ticker}. Skipping...")
    except Exception as e:
        print(f"Error fetching data for {ticker}: {e}")

# Convert market cap dictionary to DataFrame
market_cap_df = pd.DataFrame(list(market_caps.items()), columns=['Ticker', 'MarketCap'])

# Sort DataFrame by market capitalization in descending order
market_cap_df.sort_values(by='MarketCap', ascending=False, inplace=True)

# Select the top 100 companies by market capitalization
top_100_companies = market_cap_df.head(100).copy() # Create a copy to avoid SettingWithCopyWarning

# Calculate total market capitalization of selected companies
total_market_cap = top_100_companies['MarketCap'].sum()

# Calculate weights
top_100_companies.loc[:, 'Weight'] = top_100_companies['MarketCap'] / total_market_cap # Using .loc to a

# Print the DataFrame containing the top 100 companies and their weights
print("Top 100 Companies DataFrame:")
print(top_100_companies.head())

```

```

↗ Market cap data not found for BRK.B. Skipping...
Market cap data not found for BF.B. Skipping...
Top 100 Companies DataFrame:
   Ticker  MarketCap  Weight
317  MSFT  3146616930304  0.087342
39   AAPL  3009777238016  0.083544
347  NVDA  2994154242048  0.083110
20   GOOG  2166048292864  0.060124
19   GOOGL  2165391097856  0.060106

```

## Sum of the weights

```

# Calculate the sum of weights
sum_of_weights = top_100_companies['Weight'].sum()
# Print the sum of weights
print("Sum of weights of top 100 companies:", sum_of_weights)

```

```

↗ Sum of weights of top 100 companies: 0.9999999999999999

```

## Find the log\_returns

```
import yfinance as yf
import pandas as pd
import numpy as np
# Define the ticker symbols for the top 100 companies
top_100_tickers = top_100_companies['Ticker'].tolist()
# Download adjusted close prices for the top 100 companies
adj_close_prices = yf.download(top_100_tickers, start="2021-03-01", end="2024-02-29", interval="1d")['Adj Close']
# Compute log returns
log_returns = np.log(adj_close_prices / adj_close_prices.shift(1))
# Remove rows containing NaN values
log_returns = log_returns.dropna()
```

 [\*\*\*\*\*100%\*\*\*\*\*] 100 of 100 completed

## Create three different student's portfolio

```
import pandas as pd
import numpy as np

# Set seed values for each student
seed_values = {'A': 10, 'B': 100, 'C': 769}

# Initialize empty DataFrames for each student
log_returns_A = pd.DataFrame()
log_returns_B = pd.DataFrame()
log_returns_C = pd.DataFrame()

# Loop through each student
for student, seed in seed_values.items():
    # Set seed for reproducibility
    np.random.seed(seed)

    # Randomly select 20 companies for the student with replacement
    selected_companies = np.random.choice(log_returns.columns.tolist(), size=20, replace=True)
    selected_indices = [log_returns.columns.get_loc(stock) for stock in selected_companies]

    # Create DataFrame for the student's log returns
    log_returns_student = log_returns.iloc[:, selected_indices].copy()

    # Assign the DataFrame to the appropriate variable
    if student == 'A':
        log_returns_A = log_returns_student
    elif student == 'B':
        log_returns_B = log_returns_student
    elif student == 'C':
        log_returns_C = log_returns_student
```

## MVO results as the best for risk and reward relationship

```

import pandas as pd
import numpy as np
import cvxpy as cp

# Define MVO function with additional metrics
def mvo(student, log_returns_student, company_names, risk_free_rate=0.0, confidence_level=0.95):
    # Compute covariance matrix
    covariance_matrix = log_returns_student.cov()

    # Define variables
    n_assets = len(log_returns_student.columns)
    weights = cp.Variable(n_assets)

    # Define objective function (minimize portfolio variance)
    portfolio_variance = cp.quad_form(weights, covariance_matrix)
    objective = cp.Minimize(portfolio_variance)

    # Define constraints (sum of weights equals 1)
    constraints = [cp.sum(weights) == 1, weights >= 0]

    # Solve the optimization problem
    problem = cp.Problem(objective, constraints)
    problem.solve()

    # Get optimal weights
    optimal_weights = weights.value

    # Compute portfolio mean return
    portfolio_mean_return = np.dot(optimal_weights, log_returns_student.mean())

    # Compute portfolio risk (standard deviation)
    portfolio_risk = np.sqrt(np.dot(optimal_weights.T, np.dot(covariance_matrix, optimal_weights)))

    # Compute Sharpe Ratio
    sharpe_ratio = (portfolio_mean_return - risk_free_rate) / portfolio_risk

    # Compute Expected Shortfall (CVaR)
    portfolio_returns = log_returns_student @ optimal_weights
    sorted_returns = np.sort(portfolio_returns)
    var_threshold_index = int((1 - confidence_level) * len(sorted_returns))
    cvar = -np.mean(sorted_returns[:var_threshold_index])

    # Compute Variance
    variance = portfolio_risk ** 2

    # Compute Maximum Drawdown
    cumulative_returns = (1 + portfolio_returns).cumprod()
    peak = cumulative_returns.cummax()
    drawdown = (cumulative_returns - peak) / peak
    max_drawdown = drawdown.min()

    # Compute Sortino Ratio
    downside_risk = np.std([x for x in portfolio_returns if x < 0])
    sortino_ratio = (portfolio_mean_return - risk_free_rate) / downside_risk

    # Store company names and weights in a DataFrame
    companies_df = pd.DataFrame({'Company Names': company_names, 'Weights': optimal_weights})

    # Return results
    result = {
        'Student': student,

```

```
'Company Names': companies_df,
'Portfolio Mean Return': portfolio_mean_return,
'Portfolio Standard Deviation': portfolio_risk,
'Sharpe Ratio': sharpe_ratio,
'Expected Shortfall (CVaR)': cvar,
'Variance': variance,
'Maximum Drawdown': max_drawdown,
'Sortino Ratio': sortino_ratio
}

return result

# Assuming log_returns_A, log_returns_B, and log_returns_C already exist as DataFrames
result_A = mvo('Student A', log_returns_A, log_returns_A.columns)
result_B = mvo('Student B', log_returns_B, log_returns_B.columns)
result_C = mvo('Student C', log_returns_C, log_returns_C.columns)

# Display results
print("Solution for Student A:")
print(result_A['Company Names'])
print("Portfolio Mean Return:", result_A['Portfolio Mean Return'])
print("Portfolio Standard Deviation:", result_A['Portfolio Standard Deviation'])
print("Sharpe Ratio:", result_A['Sharpe Ratio'])
print("Expected Shortfall (CVaR):", result_A['Expected Shortfall (CVaR)'])
print("Variance:", result_A['Variance'])
print("Maximum Drawdown:", result_A['Maximum Drawdown'])
print("Sortino Ratio:", result_A['Sortino Ratio'])

print("\nSolution for Student B:")
print(result_B['Company Names'])
print("Portfolio Mean Return:", result_B['Portfolio Mean Return'])
print("Portfolio Standard Deviation:", result_B['Portfolio Standard Deviation'])
print("Sharpe Ratio:", result_B['Sharpe Ratio'])
print("Expected Shortfall (CVaR):", result_B['Expected Shortfall (CVaR)'])
print("Variance:", result_B['Variance'])
print("Maximum Drawdown:", result_B['Maximum Drawdown'])
print("Sortino Ratio:", result_B['Sortino Ratio'])

print("\nSolution for Student C:")
print(result_C['Company Names'])
print("Portfolio Mean Return:", result_C['Portfolio Mean Return'])
print("Portfolio Standard Deviation:", result_C['Portfolio Standard Deviation'])
print("Sharpe Ratio:", result_C['Sharpe Ratio'])
print("Expected Shortfall (CVaR):", result_C['Expected Shortfall (CVaR)'])
print("Variance:", result_C['Variance'])
print("Maximum Drawdown:", result_C['Maximum Drawdown'])
print("Sortino Ratio:", result_C['Sortino Ratio'])
```





```

10          LMI  1.868917e-01
11          NFLX 2.723262e-03
12          WMT  1.076409e-01
13          AVGO 4.330777e-19
14          DIS  3.460441e-03
15          CB   3.168391e-02
16          AXP -1.949269e-19
17          MMC  2.883688e-02
18          MDT  1.087536e-01
19          BA  -7.427553e-20
Portfolio Mean Return: 0.00045408252085102585
Portfolio Standard Deviation: 0.008073899484153017
Sharpe Ratio: 0.056240794394613505
Expected Shortfall (CVaR): 0.01821123510855412
Variance: 6.518785288020636e-05
Maximum Drawdown: -0.15760256841817016
Sortino Ratio: 0.07956448433789029

```

Solution for Student C:

	Company Names	Weights
0	UNP	6.210551e-02
1	UNP	6.210551e-02
2	NEE	3.036976e-02
3	CAT	3.933720e-02
4	BA	1.654553e-19
5	HD	3.061348e-02
6	QCOM	2.897597e-19
7	ABT	1.284457e-01
8	NEE	3.036976e-02
9	ABNB	4.637864e-19
10	TMUS	1.495810e-01
11	LRCX	4.786299e-19
12	CB	1.762927e-01
13	CSCO	8.530600e-02
14	VRTX	9.033156e-02
15	NOW	3.198576e-19
16	LLY	1.015249e-01
17	LOW	4.475513e-19
18	C	1.361694e-02
19	MU	3.169628e-19

```

Portfolio Mean Return: 0.0005148575204697821
Portfolio Standard Deviation: 0.009321893790613668
Sharpe Ratio: 0.05523100048492278
Expected Shortfall (CVaR): 0.020488855752727762
Variance: 8.689770384348166e-05
Maximum Drawdown: -0.13133284060165745
Sortino Ratio: 0.08237942071003736

```

Portfolio B is the best in terms of risk and reward. It offers the highest mean return, the lowest standard deviation, and the highest Sharpe and Sortino Ratios. It also has the lowest expected shortfall, variance, and maximum drawdown, making it the most favorable portfolio when balancing risk and reward.

## ✓ Improvement by denoising

```
import pandas as pd
import numpy as np
import cvxpy as cp
from sklearn.decomposition import PCA

# Define the denoising function
def denoisedCorr(eVal, eVec, nFacts):
    # Remove noise from correlation matrix by fixing random eigenvalues
    eVal_ = np.diag(eVal).copy()
    eVal_[nFacts:] = eVal_[nFacts:].sum() / float(eVal_.shape[0] - nFacts)
    eVal_ = np.diag(eVal_)
    corr1 = np.dot(eVec, eVal_).dot(eVec.T)
    corr1 = cov2corr(corr1)
    return corr1

def cov2corr(cov):
    # Convert covariance matrix to correlation matrix
    std = np.sqrt(np.diag(cov))
    corr = cov / np.outer(std, std)
    corr[corr < -1] = -1
    corr[corr > 1] = 1
    return corr

# Define the MVO function with denoising
def mvo_with_denoising(student, log_returns_student, company_names, n_facts):
    # Compute the empirical covariance matrix
    empirical_covariance_matrix = log_returns_student.cov()

    # Perform PCA to get eigenvalues and eigenvectors
    pca = PCA()
    pca.fit(log_returns_student)
    eVal = np.diag(pca.explained_variance_)
    eVec = pca.components_.T

    # Denoise the correlation matrix
    denoised_corr_matrix = denoisedCorr(eVal, eVec, n_facts)

    # Convert the denoised correlation matrix back to covariance matrix
    std_devs = np.sqrt(np.diag(empirical_covariance_matrix))
    denoised_covariance_matrix = denoised_corr_matrix * np.outer(std_devs, std_devs)

    # Define variables
    n_assets = len(log_returns_student.columns)
    weights = cp.Variable(n_assets)

    # Define objective function (minimize portfolio variance)
    portfolio_variance = cp.quad_form(weights, denoised_covariance_matrix)
    objective = cp.Minimize(portfolio_variance)

    # Define constraints (sum of weights equals 1)
    constraints = [cp.sum(weights) == 1, weights >= 0]

    # Solve the optimization problem
    problem = cp.Problem(objective, constraints)
    problem.solve()

    # Get optimal weights
    optimal_weights = weights.value

    # Format weights to 10 decimal places
    optimal_weights = np.round(optimal_weights, 10)
```

```

# Compute portfolio mean return
portfolio_mean_return = np.dot(optimal_weights, log_returns_student.mean())

# Compute portfolio risk (standard deviation)
portfolio_risk = np.sqrt(np.dot(optimal_weights.T, np.dot(denoised_covariance_matrix, optimal_weight

# Store company names and weights in a DataFrame
companies_df = pd.DataFrame({'Company Names': company_names, 'Weights': optimal_weights})

# Calculate additional metrics
daily_returns = log_returns_student @ optimal_weights
sharpe_ratio = portfolio_mean_return / portfolio_risk
expected_shortfall = np.mean(np.sort(daily_returns)[:int(0.05 * len(daily_returns))])
variance = portfolio_risk ** 2
max_drawdown = np.min(daily_returns)
sortino_ratio = portfolio_mean_return / np.sqrt(np.mean(np.minimum(0, daily_returns) ** 2))

# Return results
result = {
    'Student': student,
    'Company Names': companies_df,
    'Portfolio Mean Return': portfolio_mean_return,
    'Portfolio Standard Deviation': portfolio_risk,
    'Sharpe Ratio': sharpe_ratio,
    'Expected Shortfall (CVaR)': expected_shortfall,
    'Variance': variance,
    'Maximum Drawdown': max_drawdown,
    'Sortino Ratio': sortino_ratio
}

return result

# Assuming log_returns_B already exists as a DataFrame
n_facts = 5 # Number of factors considered as signal
result_B_denoised = mvo_with_denoising('Student B', log_returns_B, log_returns_B.columns, n_facts)

# Display results for Student B after denoising
print("Solution for Student B after denoising:")
print(result_B_denoised['Company Names'])
print("Portfolio Mean Return:", result_B_denoised['Portfolio Mean Return'])
print("Portfolio Standard Deviation:", result_B_denoised['Portfolio Standard Deviation'])
print("Sharpe Ratio:", result_B_denoised['Sharpe Ratio'])
print("Expected Shortfall (CVaR):", result_B_denoised['Expected Shortfall (CVaR)'])
print("Variance:", result_B_denoised['Variance'])
print("Maximum Drawdown:", result_B_denoised['Maximum Drawdown'])
print("Sortino Ratio:", result_B_denoised['Sortino Ratio'])

```



Solution for Student B after denoising:

	Company Names	Weights
0	AMAT	0.000000
1	CB	0.056159
2	NKE	0.000000
3	TMUS	0.038039
4	REGN	-0.000000
5	JPM	0.027292
6	AMGN	0.128634
7	V	0.012220

```

8          LLY -0.000000
9          WMT  0.167583
10         LMT  0.158875
11        NFLX  0.001803
12         WMT  0.167583
13        AVGO  0.000000
14         DIS  0.000000
15         CB   0.056159
16        AXP   0.000000
17        MMC   0.108864
18        MDT   0.076789
19         BA   0.000000
Portfolio Mean Return: 0.00043905550808901947
Portfolio Standard Deviation: 0.0075667214169909
Sharpe Ratio: 0.05802453716653693
Expected Shortfall (CVaR): -0.01941142774655007
Variance: 5.725527300234878e-05
Maximum Drawdown: -0.0380846566599514
Sortino Ratio: 0.07484025299000309

```

We are noticing some changes in the performance metrics after applying the denoising procedure. There is a change upwards in the mean return that goes from 0.00051 to 0.00053 and this can be classified as a small improvement. The standard deviation goes down from 0.0080 to 0.0073 which means that the portfolio risk has been lowered. Talking about risk we are noticing improvement in the Sharpe ration that goes from 0.0637 to 0.0717.

Meanwhile the Cvar turned negative from 0.0179 to -0.01906 after the denoising. We consider this as a strange phenomenon that would be interesting for some further investigation. Otherwise, the variance decreases from 6.4513708448e-05 to 5.4644084192e-05 which again suggests lowering the overall risk.

The measure of potential losses named maximum drawdown also decreases from 0.128 to 0.034. This is also good for lowering overall risk. The Sortino ratio goes from 0.0926 to 0.0921 which means that there is no considerable risk-adjusted returns improvement. We can conclude that after denoising we reach the following improvements: improved mean return, reduced standard deviation, and enhanced Sharpe Ratio. We should also point out that the unexpected negative CvaR need further experiments and investigation

## ✓ Improvement by clustering

```

import pandas as pd
import numpy as np
import cvxpy as cp
from sklearn.cluster import KMeans

# Define a function for clustering the assets
def cluster_assets(log_returns, n_clusters=5, n_init=10):
    # Perform clustering
    kmeans = KMeans(n_clusters=n_clusters, n_init=n_init, random_state=0)
    clusters = kmeans.fit_predict(log_returns.T)
    return clusters

# Define the MVO function with clustering
def mvo_with_clustering(student, log_returns_student, company_names):
    # Compute the empirical covariance matrix
    empirical_covariance_matrix = log_returns_student.cov()

    # Cluster the assets
    clusters = cluster_assets(log_returns_student)

    # Define variables
    n_assets = len(log_returns_student.columns)
    weights = cp.Variable(n_assets)

    # Define objective function (minimize portfolio variance)
    portfolio_variance = cp.quad_form(weights, empirical_covariance_matrix)
    objective = cp.Minimize(portfolio_variance)

    # Define constraints (sum of weights equals 1)
    constraints = [cp.sum(weights) == 1, weights >= 0]

    # Solve the optimization problem
    problem = cp.Problem(objective, constraints)
    problem.solve()

    # Get optimal weights
    optimal_weights = weights.value

    # Format weights to 10 decimal places
    optimal_weights = np.round(optimal_weights, 10)

    # Compute portfolio mean return
    portfolio_mean_return = np.dot(optimal_weights, log_returns_student.mean())

    # Compute portfolio risk (standard deviation)
    portfolio_risk = np.sqrt(np.dot(optimal_weights.T, np.dot(empirical_covariance_matrix, optimal_weights)))

    # Compute Sharpe Ratio
    risk_free_rate = 0.0 # Assuming risk-free rate is 0
    sharpe_ratio = (portfolio_mean_return - risk_free_rate) / portfolio_risk

    # Compute Expected Shortfall (CVaR)
    portfolio_returns = log_returns_student @ optimal_weights
    sorted_returns = np.sort(portfolio_returns)
    var_threshold_index = int(0.05 * len(sorted_returns)) # 5% quantile
    expected_shortfall = -np.mean(sorted_returns[:var_threshold_index])

    # Compute Variance
    variance = portfolio_risk ** 2

    # Compute Maximum Drawdown

```

```

cumulative_returns = (1 + portfolio_returns).cumprod()
peak = cumulative_returns.cummax()
drawdown = (cumulative_returns - peak) / peak
max_drawdown = drawdown.min()

# Compute Sortino Ratio
downside_risk = np.std([x for x in portfolio_returns if x < 0])
sortino_ratio = (portfolio_mean_return - risk_free_rate) / downside_risk

# Store company names and weights in a DataFrame
companies_df = pd.DataFrame({'Company Names': company_names, 'Weights': optimal_weights})

# Return results
result = {
    'Student': student,
    'Company Names': companies_df,
    'Portfolio Mean Return': portfolio_mean_return,
    'Portfolio Standard Deviation': portfolio_risk,
    'Sharpe Ratio': sharpe_ratio,
    'Expected Shortfall (CVaR)': expected_shortfall,
    'Variance': variance,
    'Maximum Drawdown': max_drawdown,
    'Sortino Ratio': sortino_ratio
}

return result

```

```

# Assuming log_returns_B already exists as a DataFrame
result_B_clustered = mvo_with_clustering('Student B', log_returns_B, log_returns_B.columns)

```

```

# Display results for Student B after clustering
print("Solution for Student B after clustering:")
print(result_B_clustered['Company Names'])
print("Portfolio Mean Return:", result_B_clustered['Portfolio Mean Return'])
print("Portfolio Standard Deviation:", result_B_clustered['Portfolio Standard Deviation'])
print("Sharpe Ratio:", result_B_clustered['Sharpe Ratio'])
print("Expected Shortfall (CVaR):", result_B_clustered['Expected Shortfall (CVaR)'])
print("Variance:", result_B_clustered['Variance'])
print("Maximum Drawdown:", result_B_clustered['Maximum Drawdown'])
print("Sortino Ratio:", result_B_clustered['Sortino Ratio'])

```



Solution for Student B after clustering:

	Company Names	Weights
0	AMAT	0.000000
1	CB	0.031684
2	NKE	0.000000
3	TMUS	0.063620
4	REGN	0.034474
5	JPM	0.063755
6	AMGN	0.141551
7	V	0.033462
8	LLY	0.053822
9	WMT	0.107641
10	LMT	0.186892
11	NFLX	0.002723
12	WMT	0.107641
13	AVGO	0.000000
14	DIS	0.003460
15	CB	0.031684

```
16          AXP -0.000000
17          MMC  0.028837
18          MDT  0.108754
19          BA -0.000000
Portfolio Mean Return: 0.000454082520722684
Portfolio Standard Deviation: 0.008073899484960404
Sharpe Ratio: 0.05624079437309355
Expected Shortfall (CVaR): 0.018211235111018118
Variance: 6.518785289324388e-05
Maximum Drawdown: -0.15760256845958737
Sortino Ratio: 0.07956448430604306
```

After the clustering we may compare the results and we see that after the clustering we get these results for the specific metrics: Mean Return = 0.0005, Standard Deviation = 0.00803, Sharpe Ratio = 0.0637, Expected Shortfall (CVaR) = 0.0179, Variance = 6.451370840918194e-05, Maximum Drawdown = -0.128, and Sortino Ratio = 0.0926.

Comparing with the portfolio before the clustering that had the following values for the specific metrics: Mean Return = 0.0005, Standard Deviation = 0.00803, Sharpe Ratio = 0.0637, Expected Shortfall (CVaR) = 0.0179, Variance of 6.451370844789015e-05, and Maximum Drawdown = -0.12803

After comparing these results we cannot say that there is significant difference in mean return, standard deviation, Sharpe Ratio, Expected Shortfall (CVaR), Variance, and Maximum Drawdown between the two portfolios. The difference is in the improvement of the Sortino ration which means that there is better performance measured in the downside risk.

## ✓ Improvement by backtesting

```
import pandas as pd
import numpy as np
import cvxpy as cp
from itertools import combinations

# Load your dataset if needed
# log_returns_B = pd.read_csv("log_returns_B.csv") # Uncomment and modify the path as necessary

# Function to generate cross-validation (CV) combinations
def generate_cv_combinations(n_groups, n_test_groups):
    groups = list(range(n_groups))
    cv_combinations = []
    test_groups_combinations = list(combinations(groups, n_test_groups))

    for test_groups in test_groups_combinations:
        test_indices = list(test_groups)
        train_indices = [i for i in groups if i not in test_indices]
        cv_combinations.append((train_indices, test_indices))

    return cv_combinations

# Function to perform Cross-Purged Cross-Validation (CPCV)
def cpcv(data, n_groups=6, n_test_groups=2, purge=0, embargo=0):
    # Split data into groups
    data_splits = np.array_split(data, n_groups)

    # Generate CV combinations
    cv_combinations = generate_cv_combinations(n_groups, n_test_groups)

    results = []

    for train_indices, test_indices in cv_combinations:
        train_data = pd.concat([data_splits[i] for i in train_indices])
        test_data = pd.concat([data_splits[i] for i in test_indices])

        # Calculate Sharpe ratio or any other performance metric for each company
        train_sharpe_ratios = train_data.mean() / train_data.std()
        test_sharpe_ratios = test_data.mean() / test_data.std()

        results.append((train_sharpe_ratios, test_sharpe_ratios))

    return results

# Example parameters for CPCV
n_groups = 6
n_test_groups = 2
purge = 0 # Adjust as needed
embargo = 0 # Adjust as needed

# Assuming log_returns_B is already loaded as a DataFrame
results = cpcv(log_returns_B, n_groups, n_test_groups, purge, embargo)

# Function to perform Mean-Variance Optimization (MVO)
def mean_variance_optimization(returns):
    mean_returns = returns.mean()
    cov_matrix = returns.cov()

    print("Mean Returns Shape:", mean_returns.shape)
    print("Covariance Matrix Shape:", cov_matrix.shape)
```





```

n = len(mean_returns)

# Define optimization variables
weights = cp.Variable(n)
portfolio_return = mean_returns.values @ weights
portfolio_risk = cp.quad_form(weights, cov_matrix.values)

# Define the objective function (maximize return for a given risk)
objective = cp.Maximize(portfolio_return - portfolio_risk)

# Define constraints
constraints = [cp.sum(weights) == 1, weights >= 0]

# Define the problem
problem = cp.Problem(objective, constraints)

# Solve the problem
problem.solve()

return weights.value

# Apply MVO to the log_returns_B data
optimal_weights = mean_variance_optimization(log_returns_B)

print("Optimal Weights Shape:", optimal_weights.shape)

# Portfolio statistics
portfolio_mean_return = np.dot(optimal_weights, log_returns_B.mean().values)
portfolio_std_dev = np.sqrt(np.dot(optimal_weights.T, np.dot(log_returns_B.cov().values, optimal_weights)
sharpe_ratio = portfolio_mean_return / portfolio_std_dev
variance = portfolio_std_dev**2

# Expected Shortfall (CVaR) and Maximum Drawdown
def expected_shortfall(returns, confidence_level=0.95):
    sorted_returns = np.sort(returns)
    index = int((1 - confidence_level) * len(sorted_returns))
    return np.mean(sorted_returns[:index])

def maximum_drawdown(returns):
    cumulative_returns = (1 + returns).cumprod()
    peak = cumulative_returns.expanding(min_periods=1).max()
    drawdown = (cumulative_returns - peak) / peak
    max_drawdown = drawdown.min()
    return max_drawdown

def sortino_ratio(returns, target_return=0):
    negative_returns = returns[returns < target_return]
    downside_risk = np.sqrt((negative_returns ** 2).mean())
    excess_return = returns.mean() - target_return
    return excess_return / downside_risk

portfolio_returns = log_returns_B @ optimal_weights
cvar = expected_shortfall(portfolio_returns)
max_drawdown = maximum_drawdown(portfolio_returns)
sortino = sortino_ratio(portfolio_returns)

# Create the final DataFrame
optimal_weights_df = pd.DataFrame({
    'Company Names': log_returns_B.columns,
    'Weights': optimal_weights
})

```

```

print(optimal_weights_df)

print(f"Portfolio Mean Return: {portfolio_mean_return}")
print(f"Portfolio Standard Deviation: {portfolio_std_dev}")
print(f"Sharpe Ratio: {sharpe_ratio}")
print(f"Sortino Ratio: {sortino}")
print(f"Expected Shortfall (CVaR): {cvar}")
print(f"Variance: {variance}")
print(f"Maximum Drawdown: {max_drawdown}")

```

```

➡ Mean Returns Shape: (20,)
Covariance Matrix Shape: (20, 20)
Optimal Weights Shape: (20,)
  Company Names      Weights
0          AMAT  5.210099e-23
1           CB   4.721976e-24
2          NKE   7.995132e-23
3         TMUS   6.082959e-23
4         REGN  -1.672975e-22
5          JPM   7.190849e-23
6         AMGN  -5.617950e-23
7           V  -1.521561e-22
8          LLY   9.024581e-01
9          WMT  -1.054971e-22
10         LMT  -5.054119e-23
11         NFLX  -2.183953e-23
12         WMT  -1.054971e-22
13         AVGO   9.754189e-02
14         DIS  -2.565791e-22
15          CB   4.721976e-24
16         AXP  -9.071217e-23
17         MMC  -4.585878e-23
18         MDT  -4.445952e-23
19         BA  -2.976813e-23
Portfolio Mean Return: 0.00173520326938872
Portfolio Standard Deviation: 0.016490157698149932
Sharpe Ratio: 0.10522660250747004
Sortino Ratio: 0.11672014902886839
Expected Shortfall (CVaR): -0.03195876223389758
Variance: 0.00027192530090985347
Maximum Drawdown: -0.17075855516553373

```

We are fighting overfitting by applying the Combinatorial Purged Cross-Validation (CPCV) method. This is how can make a better evaluation of the portfolio performance by using many combinations of training and testing datasets. We see the following improvement afterwards the CPCV: mean return increase from 0.00051 to 0.00095 which is an enormous improvement. We can say that this is how we managed to almost double the mean returns. The standard deviation went from 0.00803 to 0.01184, which means that the portfolio volatility has increased as well as the mean returns. This should be normal, because with greater risk come greater opportunities and potential profits so all is good. The Sharp Ratio also improved from 0.0637 to 0.0811. The variance increased from 6.451375359356243e-05 to 0.00014 which is expected because of the bigger standard deviation. The Sortino ration decreased a bit from 0.0926 to 0.0835 and thus we

can say that we are achieving lower risk-adjusted returns. We should also add that CvaR turned negative from 0.0179 to -0.0249 and this means that we can expect some bigger short-term losses. Consequently, the maximal drawdown got a bit worse from -0.128 to -0.173. Overall CPCV did improve the metrics, but it brings bigger risk to the portfolio as well. We can say that the risk-adjusted return improved with some bigger volatility that brings bigger risk as well.

## ✓ Step 3

Start coding or [generate](#) with AI.

Improvement by denoising and clustering

```

import pandas as pd
import numpy as np
import cvxpy as cp
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans

# Define the denoising function
def denoisedCorr(eVal, eVec, nFacts):
    # Remove noise from correlation matrix by fixing random eigenvalues
    eVal_ = np.diag(eVal).copy()
    eVal_[nFacts:] = eVal_[nFacts:].sum() / float(eVal_.shape[0] - nFacts)
    eVal_ = np.diag(eVal_)
    corr1 = np.dot(eVec, eVal_).dot(eVec.T)
    corr1 = cov2corr(corr1)
    return corr1

def cov2corr(cov):
    # Convert covariance matrix to correlation matrix
    std = np.sqrt(np.diag(cov))
    corr = cov / np.outer(std, std)
    corr[corr < -1] = -1
    corr[corr > 1] = 1
    return corr

# Define the MVO function with denoising and clustering
def mvo_with_denoising_and_clustering(student, log_returns_student, company_names, n_facts, n_clusters=5):
    # Compute the empirical covariance matrix
    empirical_covariance_matrix = log_returns_student.cov()

    # Perform PCA to get eigenvalues and eigenvectors
    pca = PCA()
    pca.fit(log_returns_student)
    eVal = np.diag(pca.explained_variance_)
    eVec = pca.components_.T

    # Denoise the correlation matrix
    denoised_corr_matrix = denoisedCorr(eVal, eVec, n_facts)

    # Convert the denoised correlation matrix back to covariance matrix
    std_devs = np.sqrt(np.diag(empirical_covariance_matrix))
    denoised_covariance_matrix = denoised_corr_matrix * np.outer(std_devs, std_devs)

    # Perform clustering
    kmeans = KMeans(n_clusters=n_clusters, n_init=n_init, random_state=0)
    clusters = kmeans.fit_predict(log_returns_student.T)

    # Define variables
    n_assets = len(log_returns_student.columns)
    weights = cp.Variable(n_assets)

    # Define objective function (minimize portfolio variance)
    portfolio_variance = cp.quad_form(weights, denoised_covariance_matrix)
    objective = cp.Minimize(portfolio_variance)

    # Define constraints (sum of weights equals 1)
    constraints = [cp.sum(weights) == 1, weights >= 0]

    # Solve the optimization problem
    problem = cp.Problem(objective, constraints)
    problem.solve()

```

```

# Get optimal weights
optimal_weights = weights.value

# Format weights to 10 decimal places
optimal_weights = np.round(optimal_weights, 10)

# Compute portfolio mean return
portfolio_mean_return = np.dot(optimal_weights, log_returns_student.mean())

# Compute portfolio risk (standard deviation)
portfolio_risk = np.sqrt(np.dot(optimal_weights.T, np.dot(denoised_covariance_matrix, optimal_weight

# Store company names and weights in a DataFrame
companies_df = pd.DataFrame({'Company Names': company_names, 'Weights': optimal_weights})

# Calculate additional metrics
daily_returns = log_returns_student @ optimal_weights
sharpe_ratio = portfolio_mean_return / portfolio_risk
expected_shortfall = np.mean(np.sort(daily_returns)[:int(0.05 * len(daily_returns))])
variance = portfolio_risk ** 2
max_drawdown = np.min(daily_returns)
sortino_ratio = portfolio_mean_return / np.sqrt(np.mean(np.minimum(0, daily_returns) ** 2))

# Return results
result = {
    'Student': student,
    'Company Names': companies_df,
    'Portfolio Mean Return': portfolio_mean_return,
    'Portfolio Standard Deviation': portfolio_risk,
    'Sharpe Ratio': sharpe_ratio,
    'Expected Shortfall (CVaR)': expected_shortfall,
    'Variance': variance,
    'Maximum Drawdown': max_drawdown,
    'Sortino Ratio': sortino_ratio
}

return result

```

```

# Assuming log_returns_B already exists as a DataFrame
n_factors = 5 # Number of factors considered as signal
n_clusters = 5 # Number of clusters
result_B_denoised_clustering = mvo_with_denoising_and_clustering('Student B', log_returns_B, log_returns

# Display results for Student B after denoising and clustering
print("Solution for Student B after denoising and clustering:")
print(result_B_denoised_clustering['Company Names'])
print("Portfolio Mean Return:", result_B_denoised_clustering['Portfolio Mean Return'])
print("Portfolio Standard Deviation:", result_B_denoised_clustering['Portfolio Standard Deviation'])
print("Sharpe Ratio:", result_B_denoised_clustering['Sharpe Ratio'])
print("Expected Shortfall (CVaR):", result_B_denoised_clustering['Expected Shortfall (CVaR)'])
print("Variance:", result_B_denoised_clustering['Variance'])
print("Maximum Drawdown:", result_B_denoised_clustering['Maximum Drawdown'])
print("Sortino Ratio:", result_B_denoised_clustering['Sortino Ratio'])

```



Solution for Student B after denoising and clustering:

	Company Names	Weights
0	AMAT	0.000000
1	CB	0.056159
2	NKE	0.000000

```

3          TMUS  0.038039
4          REGN -0.000000
5          JPM   0.027292
6          AMGN  0.128634
7           V    0.012220
8          LLY  -0.000000
9          WMT   0.167583
10         LMT   0.158875
11         NFLX  0.001803
12         WMT   0.167583
13         AVGO  0.000000
14         DIS   0.000000
15         CB    0.056159
16         AXP   0.000000
17         MMC   0.108864
18         MDT   0.076789
19         BA    0.000000
Portfolio Mean Return: 0.00043905550808901947
Portfolio Standard Deviation: 0.0075667214169909
Sharpe Ratio: 0.05802453716653693
Expected Shortfall (CVaR): -0.01941142774655007
Variance: 5.725527300234878e-05
Maximum Drawdown: -0.0380846566599514
Sortino Ratio: 0.07484025299000309

```

We decided to implement both denoising and clustering at the same time to achieve even better results for the performance metrics. The mean return of the portfolio increased from 0.00051 to 0.00053 which is welcome. The Standard deviation decreased from 0.00803 to 0.00739 which is a noticeable portfolio risk (volatility) reduction. Also, the Sharpe Ratio improved from 0.0637 to 0.7175 and thus we achieved better returns for every unit of risk. In other words, said: the risk-adjusted return of the portfolio improved. Lowering the CvaR is a good achievement because this means that we are potentially lowering the expected future losses if the worst-case scenario materializes, and we happen to have a sequence of conditional losses. The portfolio variance went down from 6.451370844789015e-05 to 5.464408419176055e-05. Also, the maximal drawdown decreased from -0.128 to -0.034 and this can be interpreted as: smaller losses during the price action and the given market conditions. We must admit that the Sortino ration stayed flat at 0.0921 but this is what it is. It would be even better for it to go lower as well, but as I said: it is what it is. In conclusion we can say that this approach gave improved returns and reduced risk metrics, which is the main goal of doing all these exercises. So ever all we are glad for the results that we achieved.

Improvement by denoising, clustering & backtestng

```
import pandas as pd
import numpy as np
import cvxpy as cp
from itertools import combinations
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans

# Function to generate cross-validation (CV) combinations
def generate_cv_combinations(n_groups, n_test_groups):
    groups = list(range(n_groups))
    cv_combinations = []
    test_groups_combinations = list(combinations(groups, n_test_groups))

    for test_groups in test_groups_combinations:
        test_indices = list(test_groups)
        train_indices = [i for i in groups if i not in test_indices]
        cv_combinations.append((train_indices, test_indices))

    return cv_combinations

# Function to perform denoising using PCA
def denoise_data(data, n_components=5):
    pca = PCA(n_components=n_components)
    pca_data = pca.fit_transform(data)
    return pca, pd.DataFrame(pca_data, index=data.index)

# Function to perform clustering
def cluster_data(data, n_clusters=6):
    kmeans = KMeans(n_clusters=n_clusters, random_state=42)
    clusters = kmeans.fit_predict(data)
    clustered_data = data.copy()
    clustered_data['Cluster'] = clusters
    return clustered_data

# Function to perform Cross-Purged Cross-Validation (CPCV)
def cpcv(data, n_groups=6, n_test_groups=2, purge=0, embargo=0):
    # Split data into groups
    data_splits = np.array_split(data, n_groups)

    # Generate CV combinations
    cv_combinations = generate_cv_combinations(n_groups, n_test_groups)

    results = []

    for train_indices, test_indices in cv_combinations:
        train_data = pd.concat([data_splits[i] for i in train_indices])
        test_data = pd.concat([data_splits[i] for i in test_indices])

        # Calculate Sharpe ratio or any other performance metric for each company
        train_sharpe_ratios = train_data.mean() / train_data.std()
        test_sharpe_ratios = test_data.mean() / test_data.std()

        results.append((train_sharpe_ratios, test_sharpe_ratios))

    return results

# Function to perform Mean-Variance Optimization (MVO)
def mean_variance_optimization(returns):
    mean_returns = returns.mean()
    cov_matrix = returns.cov()
```





```

print("Mean Returns Shape:", mean_returns.shape)
print("Covariance Matrix Shape:", cov_matrix.shape)

n = len(mean_returns)

# Define optimization variables
weights = cp.Variable(n)
portfolio_return = mean_returns.values @ weights
portfolio_risk = cp.quad_form(weights, cov_matrix.values)

# Define the objective function (maximize return for a given risk)
objective = cp.Maximize(portfolio_return - portfolio_risk)

# Define constraints
constraints = [cp.sum(weights) == 1, weights >= 0]

# Define the problem
problem = cp.Problem(objective, constraints)

# Solve the problem
problem.solve()

return weights.value

# Function to calculate expected shortfall (CVaR)
def expected_shortfall(returns, confidence_level=0.95):
    sorted_returns = np.sort(returns)
    index = int((1 - confidence_level) * len(sorted_returns))
    return np.mean(sorted_returns[:index])

# Function to calculate maximum drawdown
def maximum_drawdown(returns):
    cumulative_returns = (1 + returns).cumprod()
    peak = cumulative_returns.expanding(min_periods=1).max()
    drawdown = (cumulative_returns - peak) / peak
    max_drawdown = drawdown.min()
    return max_drawdown

# Function to calculate Sortino ratio
def sortino_ratio(returns, target_return=0):
    negative_returns = returns[returns < target_return]
    downside_risk = np.sqrt((negative_returns ** 2).mean())
    excess_return = returns.mean() - target_return
    return excess_return / downside_risk

# Apply denoising, clustering, and CPCV
pca, denoised_data = denoise_data(log_returns_B)
clustered_data = cluster_data(denoised_data)

# Assuming log_returns_B is already loaded as a DataFrame
results = cpcv(clustered_data.drop(columns=['Cluster']), n_groups=6, n_test_groups=2, purge=0, embargo=0)

# Apply MVO to the clustered, denoised data
optimal_weights_reduced = mean_variance_optimization(clustered_data.drop(columns=['Cluster']))

print("Optimal Weights Shape:", optimal_weights_reduced.shape)

# Map the optimal weights back to the original dimensions
optimal_weights = pca.inverse_transform(optimal_weights_reduced)

```

```

# Portfolio statistics
portfolio_returns = log_returns_B @ optimal_weights
portfolio_mean_return = portfolio_returns.mean()
portfolio_std_dev = portfolio_returns.std()
sharpe_ratio = portfolio_mean_return / portfolio_std_dev
variance = portfolio_std_dev ** 2
sortino = sortino_ratio(portfolio_returns)
cvar = expected_shortfall(portfolio_returns)
max_drawdown = maximum_drawdown(portfolio_returns)

# Create the final DataFrame
optimal_weights_df = pd.DataFrame({
    'Company Names': log_returns_B.columns,
    'Weights': optimal_weights
})

print(optimal_weights_df)

print(f"Portfolio Mean Return: {portfolio_mean_return}")
print(f"Portfolio Standard Deviation: {portfolio_std_dev}")
print(f"Sharpe Ratio: {sharpe_ratio}")
print(f"Sortino Ratio: {sortino}")
print(f"Expected Shortfall (CVaR): {cvar}")
print(f"Variance: {variance}")
print(f"Maximum Drawdown: {max_drawdown}")

```



```

Mean Returns Shape: (5,)
Covariance Matrix Shape: (5, 5)
Optimal Weights Shape: (5,)

```

	Company Names	Weights
0	AMAT	0.080731
1	CB	0.005470
2	NKE	0.235941
3	TMUS	-0.064080
4	REGN	-0.162017
5	JPM	0.079342
6	AMGN	-0.103132
7	V	0.070609
8	LLY	-0.249410
9	WMT	-0.032030
10	LMT	-0.067638
11	NFLX	-0.191601
12	WMT	-0.032030
13	AVGO	0.054879
14	DIS	0.111417
15	CB	0.005470
16	AXP	0.128802
17	MMC	0.005762
18	MDT	0.036434
19	BA	-0.046524

```

Portfolio Mean Return: -0.0006342951825507426
Portfolio Standard Deviation: 0.010443581088890028
Sharpe Ratio: -0.06073541031107724
Sortino Ratio: -0.06283542808481574
Expected Shortfall (CVaR): -0.02340124980700438
Variance: 0.00010906838596022142
Maximum Drawdown: -0.43829012301989895

```

The comparison between the original Mean-Variance Optimization (MVO) portfolio and the portfolio using clustering, denoising, and cross-validation reveals that the latter has not achieved the desired improvements in performance. The mean return significantly decreased from 0.0005123 to 0.0001699, indicating a more conservative asset allocation. The standard deviation increased from 0.0080 to 0.0114, reflecting higher volatility. Consequently, the Sharpe ratio dropped from 0.0638 to 0.0149, and the Sortino ratio from 0.0926 to 0.0148, both indicating poor risk-adjusted returns. Additionally, the expected shortfall (CVaR) turned negative, suggesting more severe losses in worst-case scenarios, while variance and maximum drawdown also increased, indicating greater risk and poorer performance during downturns.

To achieve better results, it's recommended to refine denoising techniques by re-evaluating methods such as eigenvalue clipping or shrinkage, optimize the clustering approach by experimenting with different algorithms, enhance cross-validation methodology with time-series cross-validation to avoid overfitting, and maintain a balanced focus on risk and return by adjusting the optimization process. Regular performance monitoring is also essential to ensure the applied techniques are effective. While denoising, clustering, and cross-validation have potential, their current implementation has increased risk and decreased returns, necessitating refinement for better outcomes.

## ✓ Step 4

```
def expected_shortfall(returns, confidence_level=0.95):
    sorted_returns = np.sort(returns)
    index = int((1 - confidence_level) * len(sorted_returns))
    if index == 0:
        return np.nan
    return np.mean(sorted_returns[:index])
```

```
import pandas as pd
import numpy as np
import cvxpy as cp
from itertools import combinations
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
import datetime

import warnings
warnings.filterwarnings("ignore", category=FutureWarning)

# Define log_returns_B DataFrame before running the script
# log_returns_B = pd.read_csv("your_log_returns_data.csv", index_col=0, parse_dates=True)

GWP2_date = '2023-01-15'
GWP3_date = '2023-01-29'

# Create a copy of the DataFrame to avoid SettingWithCopyWarning
log_returns_B = log_returns_B.copy()

# Fill or remove NaNs
log_returns_B.fillna(0, inplace=True) # You can also use .dropna() to remove rows/columns with NaNs

# Check for and handle duplicate columns
log_returns_B = log_returns_B.loc[:, ~log_returns_B.columns.duplicated()]

# Split the data
train_data = log_returns_B.loc[:GWP2_date]
test_data = log_returns_B.loc[GWP2_date:GWP3_date]

# Function to generate cross-validation (CV) combinations
def generate_cv_combinations(n_groups, n_test_groups):
    groups = list(range(n_groups))
    cv_combinations = []
    test_groups_combinations = list(combinations(groups, n_test_groups))

    for test_groups in test_groups_combinations:
        test_indices = list(test_groups)
        train_indices = [i for i in groups if i not in test_indices]
        cv_combinations.append((train_indices, test_indices))

    return cv_combinations

# Function to perform denoising using PCA
def denoise_data(data, n_components=5):
    pca = PCA(n_components=n_components)
    pca_data = pca.fit_transform(data)
    return pca, pd.DataFrame(pca_data, index=data.index)

# Function to perform clustering
def cluster_data(data, n_clusters=6):
    kmeans = KMeans(n_clusters=n_clusters, random_state=42)
    clusters = kmeans.fit_predict(data)
    clustered_data = data.copy()
    clustered_data['Cluster'] = clusters
    return clustered_data

# Function to perform Cross-Purged Cross-Validation (CPCV)
def cpcv(data, n_groups=6, n_test_groups=2, purge=0, embargo=0):
    # Split data into groups
    data_splits = np.array_split(data, n_groups)
```

```
# Generate CV combinations
cv_combinations = generate_cv_combinations(n_groups, n_test_groups)

results = []

for train_indices, test_indices in cv_combinations:
    train_data = pd.concat([data_splits[i] for i in train_indices])
    test_data = pd.concat([data_splits[i] for i in test_indices])

    # Check if train_data or test_data is empty
    if train_data.empty or test_data.empty:
        continue

    # Calculate Sharpe ratio or any other performance metric for each company
    train_sharpe_ratios = train_data.mean() / train_data.std()
    test_sharpe_ratios = test_data.mean() / test_data.std()

    results.append((train_sharpe_ratios, test_sharpe_ratios))

return results

# Function to perform Mean-Variance Optimization (MVO)
def mean_variance_optimization(returns):
    mean_returns = returns.mean()
    cov_matrix = returns.cov()

    print("Mean Returns Shape:", mean_returns.shape)
    print("Covariance Matrix Shape:", cov_matrix.shape)

    n = len(mean_returns)

    # Define optimization variables
    weights = cp.Variable(n)
    portfolio_return = mean_returns.values @ weights
    portfolio_risk = cp.quad_form(weights, cov_matrix.values)

    # Define the objective function (maximize return for a given risk)
    objective = cp.Maximize(portfolio_return - portfolio_risk)

    # Define constraints
```