

FULL LEGAL NAME	LOCATION (COUNTRY)	EMAIL ADDRESS	MARK X FOR ANY NON-CONTRIBUTING MEMBER
Kwan Wei Ma	Hong Kong	kwanweima@gmail.com	
Alper Ülkü	Turkey	alperulkul1970@gmail.com	

Statement of integrity: By typing the names of all group members in the text boxes below, you confirm that the assignment submitted is original work produced by the group (excluding any non-contributing members identified with an "X" above).

Team member 1	Kwan Wei Ma
Team member 2	Alper Ülkü
Team member 3	

Use the box below to explain any attempts to reach out to a non-contributing member. Type (N/A) if all members contributed.

Note: You may be required to provide proof of your outreach to non-contributing members upon request.

Team member 3 was withdrawn by the system.

Step 1: Summary of Huo paper

The exploration versus exploitation dilemma describes the trade-off in sequential decision-making under uncertainty between acting ambitiously to acquire new knowledge and behaving conservatively to benefit from existing knowledge. This is a frequent difficulty encountered in a variety of real-world situations, such as portfolio selection, clinical trials, online advertising, and robotics. In a single-state Markov decision process (MDP), the stochastic multi-armed bandit issue offers a mathematical framework for examining this conundrum.

The classic "multi-armed bandit" problem in portfolio selection is impacted by the inclusion of risk awareness since it introduces a new goal of taking into account risk rather than just maximizing cumulative reward. This variation has a tight connection to the portfolio selection issue, where risk management is a crucial consideration. The authors of this paper provide the Risk-Aware Upper Confidence Bound (RA-UCB) algorithm to solve the risk-aware multi-armed bandit issue with application to portfolio selection. They add risk awareness into the traditional problem.

In this paper the RA-UCB algorithm can be summarized as follows:

1. Compute the upper confidence bound (UCB) for each asset i:

$$UCB_i(t) = \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(t)}{n_i(t-1)}} + \lambda \sqrt{\hat{\sigma}_i^2(t-1)}$$

where $\hat{\mu}_i(t-1)$ is the empirical mean reward of asset i up to trial t-1, $n_i(t-1)$ is the number of times asset i has been selected up to trial t-1, $\hat{\sigma}_i^2(t-1)$ is the empirical variance of asset i up to trial t-1, and λ is the risk parameter.

2. Select the asset with the highest UCB:

$$I_t = \arg \max_i UCB_i(t)$$

3. Observe the reward $R_{I_t,t}$ and update the empirical mean and variance for asset I_t :

$$\begin{aligned} \hat{\mu}_{I_t}(t) &= \frac{\hat{\mu}_{I_t}(t-1) \cdot n_{I_t}(t-1) + R_{I_t,t}}{n_{I_t}(t-1) + 1} \\ \hat{\sigma}_{I_t}^2(t) &= \frac{\hat{\sigma}_{I_t}^2(t-1) \cdot n_{I_t}(t-1) + (R_{I_t,t} - \hat{\mu}_{I_t}(t-1))^2}{n_{I_t}(t-1) + 1} \end{aligned}$$

4. Compute the portfolio weight for each asset i:

$$w_i(t) = \frac{(1-\lambda) \cdot UCB_i(t)}{\sum_{j=1}^K (1-\lambda) \cdot UCB_j(t)}$$

5. Rebalance the portfolio to the target weights $w(t) = (w_1(t), w_2(t), \dots, w_K(t))$.

The final portfolio weights $w(T)$ are then outputted.

Figure 1: RA-UCB portfolio selection pseudocode (Huo et.al)

There are a number of potential uses for the research on the risk-aware multi-armed bandit issue with application to portfolio selection in the machine learning and quantitative finance sectors. It can be used, for instance, to improve internet advertising, where marketers must strike a balance between testing out new ads and capitalizing on the most effective ones. It can also be utilized in clinical studies, where it's important for researchers to strike a balance between using the best medicines and researching novel ones. The suggested algorithm can be used to optimize portfolio selection in quantitative finance, where investors must balance the investigation of new assets with the exploitation of the best-performing assets while taking risk management into account.

Step 2: Portfolio selection pseudocode

In order to achieve the investor's objectives, such as maximizing cumulative return in respect to some risk measure, the portfolio selection problem seeks to determine the best mix of assets that should be maintained in the portfolio. We alter the conventional multi-armed bandit paradigm to simulate portfolio selection by considering a financial market with a significant collection of assets, from which the learner selects a basket of K assets to invest in a series of N trials. Figure 1 below depicts the model's pseudocode. The learner's investment strategy consists of N mappings from the newly acquired knowledge to W in total.

Input:

- δ : length of the available history for estimating mean returns
- N : trading horizon

Historical period:

1. Download prices for K assets
2. Calculate historical returns $H_{i,t}, t = 1, \dots, \delta$

Investment period:

for t in $1, \dots, N$:

1. Choose portfolio $\omega_t = (\omega_{1,t}, \dots, \omega_{K,t})^T$ where
 - $\omega_t \in W$
 - $W = \{u \in \mathbb{R}_+^K : u^T \mathbf{1} = 1\}$
 - $\mathbf{1}$ is a column vector of ones
2. Observe $R_t = (R_{1,t}, \dots, R_{K,t})^T$ where
 - $R_{i,t} = \ln(P_{i,t+1}/P_{i,t})$
3. Receive reward $\omega_t^T R_t$

Figure 2: Sequential portfolio selection pseudocode (Model 1)

Step 3: Import and Structure Data for 2008

The sample data will be based on the daily returns of 27 selected assets between September and October 2008. The following are the asset lists for 13 financial institutions: **JPM, WFC, BAC, C, GS, USB, MS, KEY, PNC, COF, AXP, PRU, SCHW**. As STU and BBT merged in 2019, downloading the data of the two assets BBT and STI from Yahoo Finance was not possible. Thus we skipped these two.

14 non-financial institutions: **KR, PFE, XOM, WMT, DAL, CSCO, HCP, EQIX, DUK, NFLX; APA, F, REGN, CMS**. As HCP was listed in 2021, data of HCP in 2008 was not available. Thus we skipped the HCP. Disregarding these 3 assets, we have a total of 27 assets thus continue our calculations with 27 assets.

We combined the data into a suitable Python time series data structure as follows. We defined a function called:

```
DownloadAssetsConvertDataframe(tickers, startTime, endTime)
```

That outputs the prices and log returns as time series data frames called **Prices0** and **Returns0**. Here are some lines from **Prices0** dataframe:

1	Prices0
2	
	JPM WFC BAC C GS USD MS KEY PNC COF ... DAL CSCO EQIX DUK NFLX GE APA F REGN CMS
2008-09-02	26.68 21.06 25.74 150.86 130.56 1.71 30.49 8.48 49.77 34.99 ... 8.13 16.48 63.43 26.49 4.41 115.83 85.02 2.71 20.58 8.23
2008-09-03	27.17 20.92 26.52 154.81 132.37 1.57 31.14 8.56 50.23 35.57 ... 8.07 16.18 62.80 26.15 4.42 116.00 85.93 2.74 21.80 8.13
2008-09-04	25.94 20.02 24.63 144.47 127.07 1.44 29.79 8.02 49.30 33.75 ... 7.94 15.46 60.73 26.41 4.27 112.47 87.97 2.64 20.38 8.14
2008-09-05	27.09 21.05 25.94 150.55 128.92 1.47 30.54 8.73 50.35 34.83 ... 7.81 15.45 60.72 26.44 4.24 113.20 89.09 2.65 19.07 8.06
2008-09-08	28.43 22.64 27.95 160.41 134.05 1.47 31.95 9.25 52.03 37.96 ... 7.61 16.22 60.73 27.30 4.31 118.11 87.77 2.73 18.90 8.31
2008-09-09	27.00 21.03 26.17 149.05 127.68 1.36 29.83 8.85 49.37 35.44 ... 7.29 15.94 58.05 27.65 4.11 114.17 80.39 2.64 18.52 8.11
2008-09-10	26.96 21.39 26.07 147.47 124.46 1.37 28.74 8.46 48.02 35.50 ... 6.98 15.95 56.94 27.30 4.11 114.05 84.24 2.68 19.43 8.19
2008-09-11	28.50 22.84 26.61 146.91 124.02 1.37 28.58 8.62 49.49 36.09 ... 7.17 16.00 59.62 27.83 4.12 114.33 87.16 2.81 20.58 8.21
2008-09-12	28.17 23.14 27.15 141.78 121.79 1.37 27.49 9.00 49.45 35.78 ... 7.21 16.28 60.28 27.95 4.21 108.61 89.86 2.95 20.71 8.34
2008-09-15	25.31 20.92 21.37 120.31 107.01 1.28 23.77 8.34 47.72 34.82 ... 7.14 15.53 56.61 27.28 4.02 99.88 83.46 2.85 20.45 8.00
2008-09-16	27.87 23.57 23.78 124.34 105.05 1.29 21.19 8.50 52.03 38.20 ... 8.81 15.82 58.67 27.56 4.07 101.75 87.02 3.04 21.64 7.67

And this ends with:

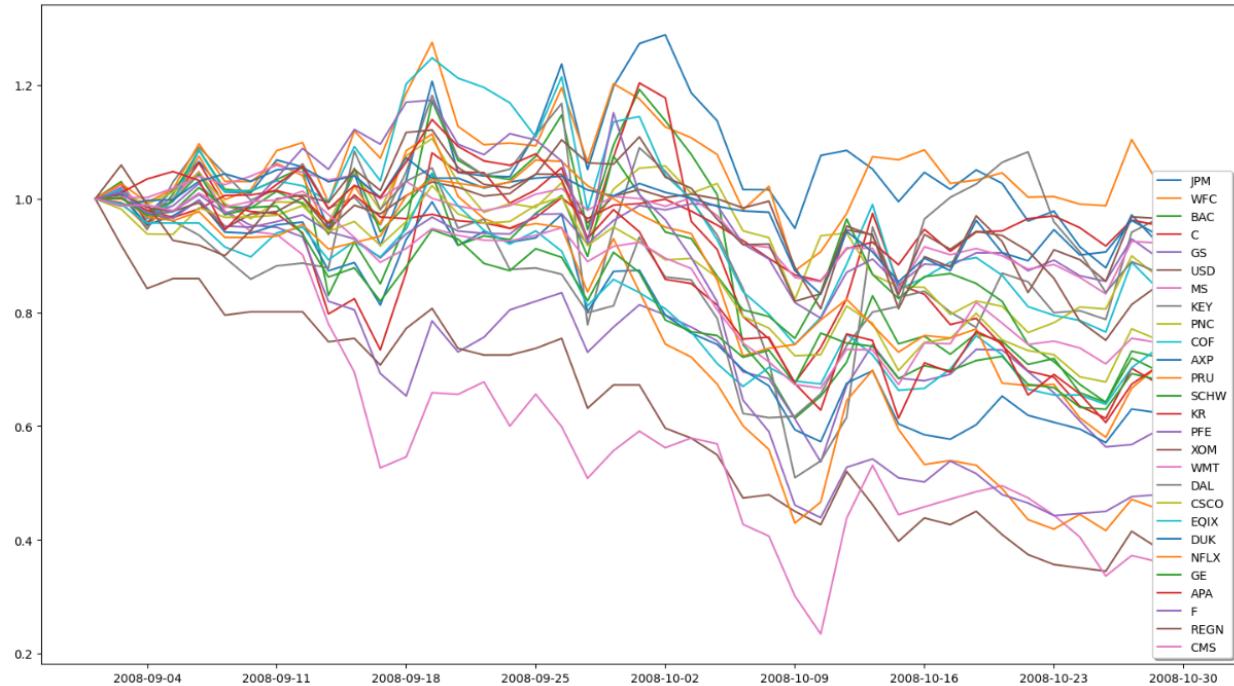


Figure 3: Price time series for 27 assets (2008 data)

Step 4: Correlation Matrix

Making use of the python code written within the defined function, `PlotCorrelations(log_returns, CORR_FILTER)`, we managed to have the cross correlation matrix together with heatmaps.

The construction of the correlation matrices is depicted below. Due to the size of the data, a heatmap has been created for easier visualization of data.

these assets, as the portfolio diversification would be suboptimal if we chose abundantly many assets from these financial institutions .

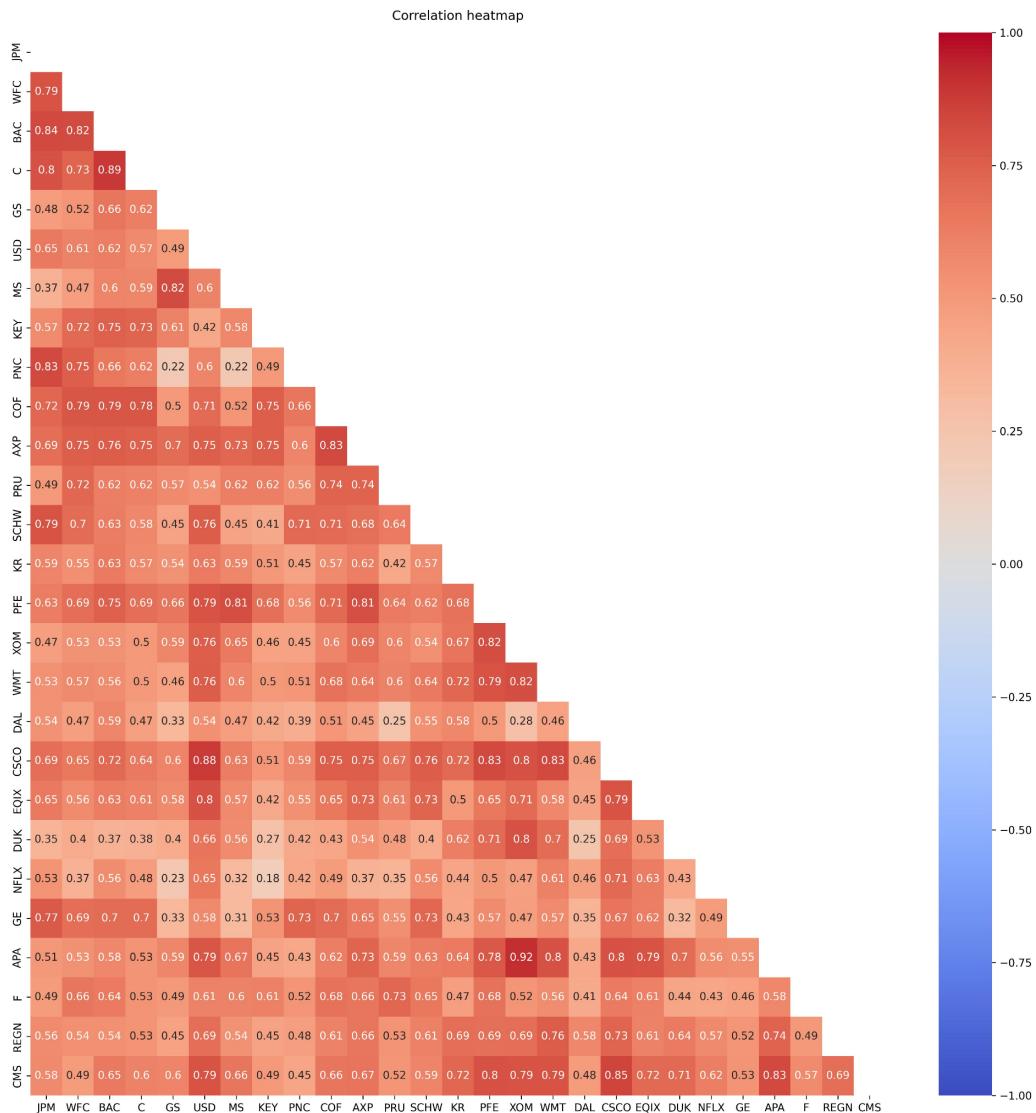


Figure 5: Triangular 27x27 correlation matrix with heatmap

2. This time we use only the lower triangular half for making assets more visible for their cross correlations.

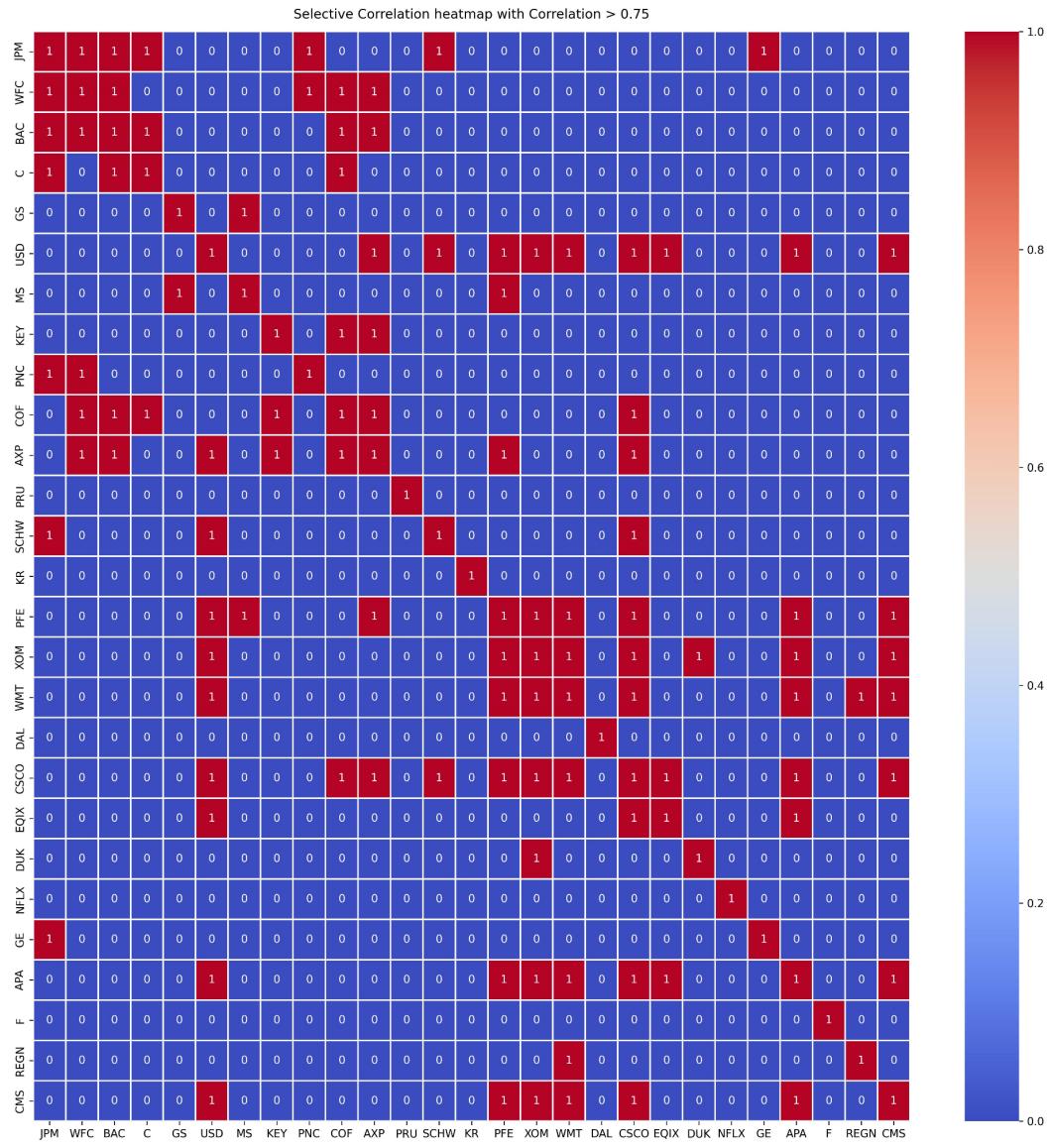


Figure 6: 27x27 correlation matrix with tagged correlations greater than 0.75

3. Here we filtered the assets for high correlation for tidying them up so that we make the discrimination for CORR_FILTER > 0.75 and filter them next time, so that these assets will not appear in the same portfolio.

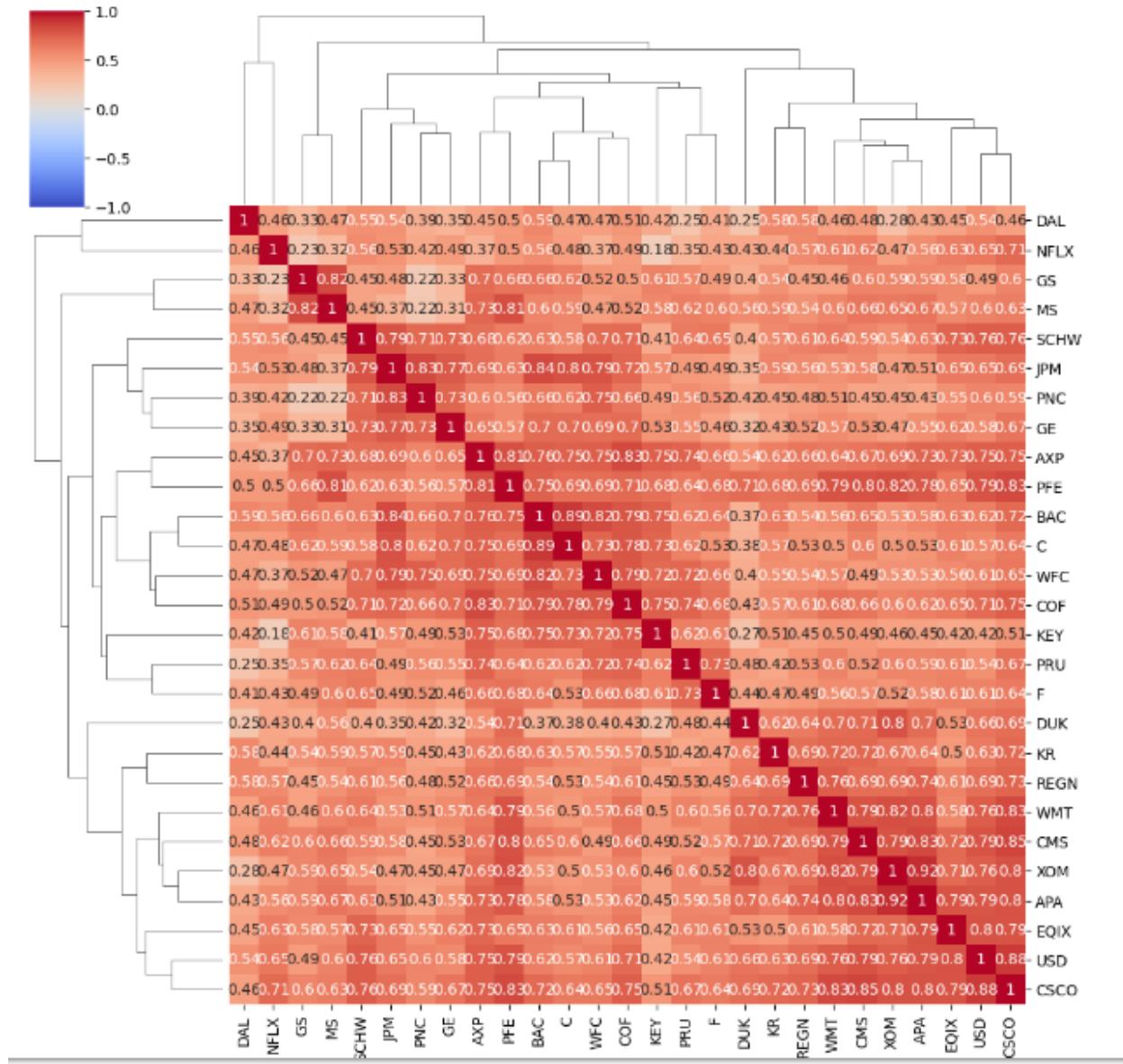


Figure 7: 27x27 correlation matrix with cluster map

4. This time by clustermap function we grouped similar correlations closer to each other. From this figure we can easily see that these groups of assets have similar highly correlated to each other.
 - a. Group-1: [AXP, PFE]
 - b. Group-2: [BAC,C, WFC, COF, KEY, PRU],
 - c. Group-3: [DUK, KR, REGN, WMT, CMS, XDM, APA],
 - d. Group-4: [SCHW, JPM, PNC, GE],
 - e. Group-5: [GS, MS]
5. Now one best asset selection from each group is required so that our portfolio can be diversified much easier and in an optimal manner.

Step 6: Upper-confidence bound (UCB) algorithm

Initialize the parameters:

NK = Number of arms (actions)

EPSILON = Epsilon value for exploration

ALPHA = Learning rate

NEPISODES = Number of episodes

HOLD = Number of days to hold an action

TMAX = Total number of time steps

UCB_weight = Weight for UCB adjustment

Choosing a random action with probability of initialized epsilon and selecting the optimal action given information about the past:

If the random generated number is less than or equal to the initialized epsilon greedy parameter, choose randomly across actions with equal probability.

If not, choose the action with the highest q-value using formula below:

$$A_t = \arg \max_a Q_t(a)$$

Updating the expected reward from choosing the given action based on formula below:

$$Q_t(a) = Q_{t-1}(a) + \frac{1}{N_t(a)} [r_t - Q_{t-1}(a)]$$

Loop the optimal selection and reward updating part over initialized time steps and episodes.

During the loop, we prioritize the exploration of actions that we are more uncertain about, we follow the next optimal choice using formula below by setting an upper-confidence bound c.

$$A_t = \arg \max_a \left\{ Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right\}$$

Step 8: Epsilon-greedy algorithm

Initialize the parameters:

NK = Number of arms (actions)

EPSILON = Epsilon value for exploration

ALPHA = Learning rate

NEPISODES = Number of episodes

HOLD = Number of days to hold an action

TMAX = Total number of time steps

Choosing a random action with probability of initialized epsilon and selecting the optimal action given information about the past:

If the random generated number is less than or equal to the initialized epsilon greedy parameter, choose randomly across actions with equal probability.

If not, choose the action with the highest q-value using formula below:

$$A_t = \arg \max_a Q_t(a)$$

Updating the expected reward from choosing the given action based on formula below:

$$Q_t(a) = Q_{t-1}(a) + \frac{1}{N_t(a)} [r_t - Q_{t-1}(a)]$$

Loop the optimal selection and reward updating part over initialized time steps and episodes.

Step 9: Results Analysis and Comparison

Results of the UCB method:

From the results, the strategy basically performs similar to an equal-weighted portfolio. However, it fails to track the most significant increases in returns across the sample period.

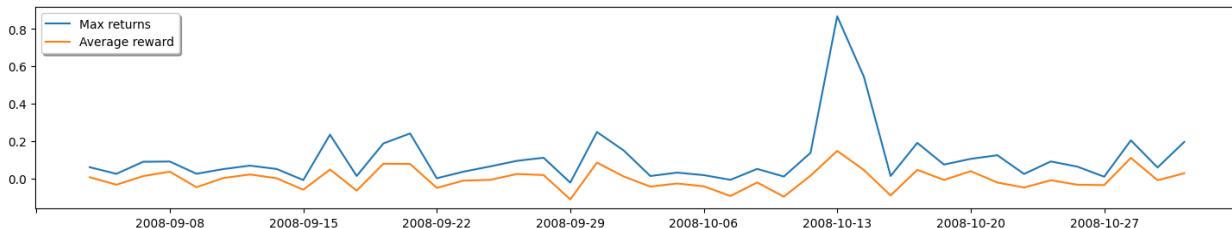


Figure 8: Average reward vs Max returns (UCB method)

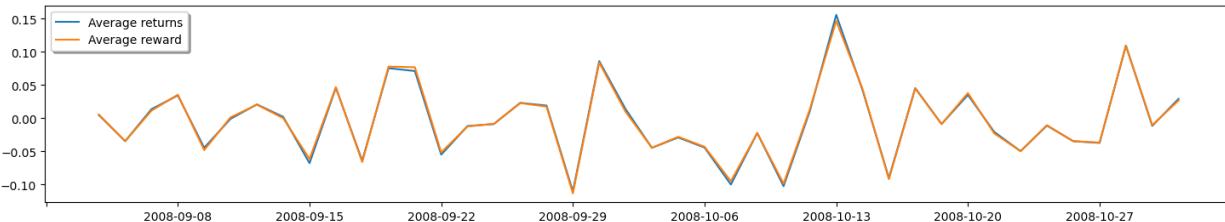


Figure 9: Average reward vs Average returns (UCB method)

We evaluate the performance of the strategy mainly using two metrics: the average frequency of optimal action selection and the average annualized return. The average frequency of optimal action selection assesses the proportion of times the algorithm selects the arm with the highest expected reward while the average annualized return measures the overall performance of the algorithms in terms of rewards generated over the investment horizon.

The average frequency of optimal action using the UCB method is 0.036 while the average annualized return from holding the Bandit portfolio is -0.6138, which is slightly better than the performance of holding the equally-weighted portfolio with the average annualized return equal to -0.6263.

Average frequency of optimal action = 0.036
 Holding the equally-weighted portfolio, average annualized return = -0.6263; average annualized standard deviation = 0.8868.
 Holding the Bandit portfolio, average annualized return = -0.6138; average annualized standard deviation = 0.8756.

The graph shows the value that would result from starting with \$1 and investing it in either the bandit approach or an equally weighted portfolio over time. Similar returns are produced by the buy-and-hold strategy of the equally weighted portfolio and the "bandit" strategy.

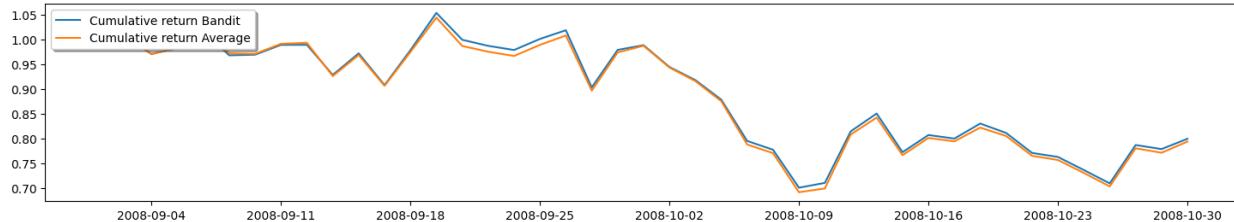


Figure 10: Cumulative return of Bandit portfolio vs equally-weighted portfolio (UCB method)

Results of Epsilon-greedy method:

From the results, the strategy basically performs similar to an equal-weighted portfolio. However, it fails to track the most significant increases in returns across the sample period.

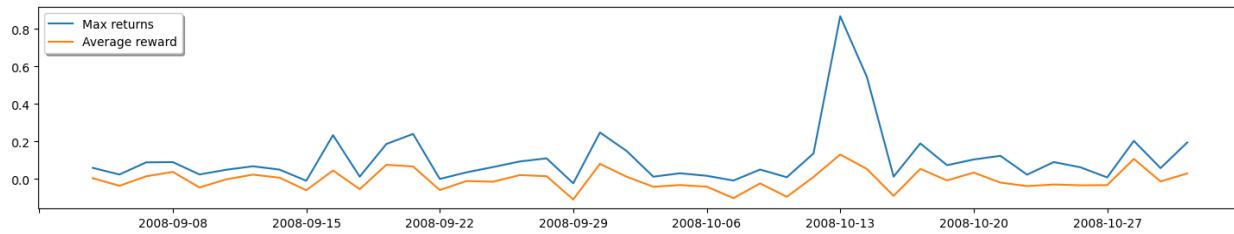


Figure 11: Average reward vs Max returns (Epsilon-greedy method)

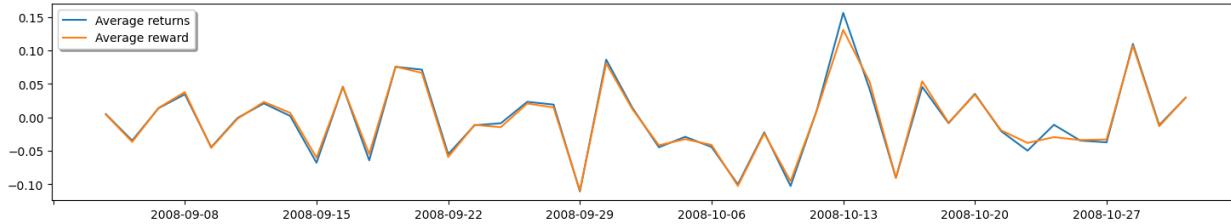


Figure 12: Average reward vs Average returns (Epsilon-greedy method)

We evaluate the performance of the strategy mainly using two metrics: the average frequency of optimal action selection and the average annualized return. The average frequency of optimal action selection assesses the proportion of times the algorithm selects the arm with the highest expected reward while the average annualized return measures the overall performance of the algorithms in terms of rewards generated over the investment horizon.

The average frequency of optimal action using Epsilon-greedy method is 0.039 while the average annualized return from holding the Bandit portfolio is -0.6491, which is slightly worse than the performance of holding the equally-weighted portfolio with the average annualized return equal to -0.6263.

Average frequency of optimal action = 0.039
 Holding the equally-weighted portfolio, average annualized return = -0.6263; average annualized standard deviation = 0.8868.
 Holding the Bandit portfolio, average annualized return = -0.6491; average annualized standard deviation = 0.8489.

The figure displays the value over time of initiating with \$1 and investing it either in the bandit strategy or the equally weighted portfolio. The ‘bandit’ strategy generates similar returns to the buy-and-hold strategy of the equally weighted portfolio.

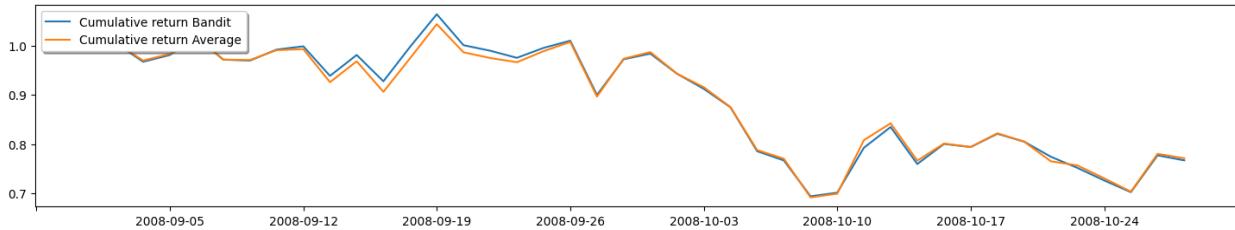


Figure 13: Cumulative return of Bandit portfolio vs equally-weighted portfolio (Epsilon-greedy method)

Comparison between results of the UCB method and results of Epsilon-greedy method:

From the results, both methods fail to track the most significant increases in returns across the sample period.

The Epsilon-greedy method shows a slightly higher average frequency of optimal action compared with the UCB method. However, it generates a lower average annualized return.

Both methods generate similar cumulative returns compared with the buy-and hold strategy of the equally weighted portfolio.

Comparison between group's result with results of the Huo paper:

The RA-UCB method is an expansion of the standard UCB algorithm that addresses the portfolio selection problem with risk management. The regular UCB algorithm balances the exploration of novel arms with the exploitation of the best-performing arms in order to maximize the cumulative reward. To do this, it calculates the UCB for each arm, chooses the one with the highest UCB, and then updates the empirical mean and variance for the chosen arm. Risk management was not specifically taken into account by our algorithm.

The RA-UCB algorithm, on the other hand, tries to maximize cumulative reward while simultaneously limiting risk. To do this, the UCB for each asset is calculated, the asset with the highest UCB is chosen, the empirical mean and variance are updated for the chosen asset, the portfolio weights are calculated for each asset, and the portfolio is then rebalanced to the target weights. The risk parameter lambda regulates the trade-off between reward and risk; a larger value of lambda indicates a choice for greater risk, while a lower value indicates a preference for reduced risk.

And we visualize the cumulative performance of each stock during the sample period.

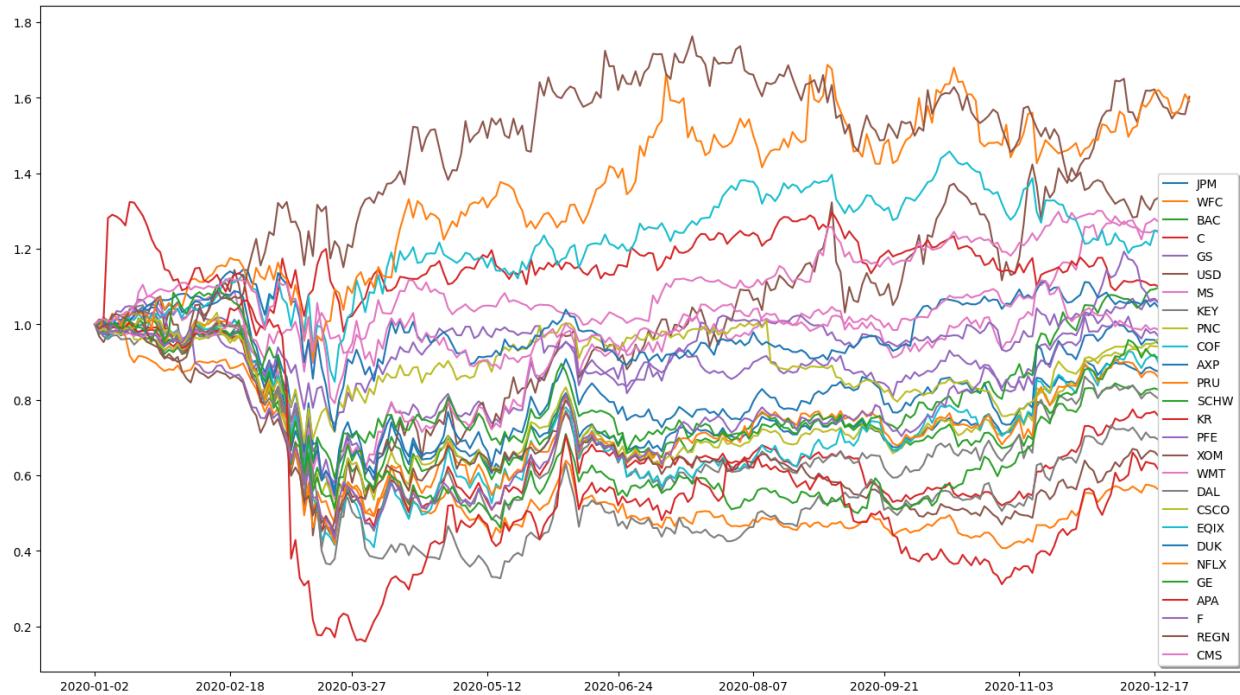


Figure 14: Price time series for 27 assets (2020 data)

Step 11: Rerun algorithm for new data (2020)

We rerun the same code in Step 6 and Step 8 for new data (2020). In order to allow more exploration, we adjust the parameters using higher epsilon and higher upper confidence bound. We also changed the holding period from 1 to 2.

Most of the results show a similar situation to the results using the data in 2008. For example, both methods fail to track the most significant increases in returns across the sample period and the Epsilon-greedy method shows a slightly higher average frequency of optimal action compared with the UCB method (0.0407 vs 0.0369).



Figure 15: Average reward vs Max returns (UCB method)

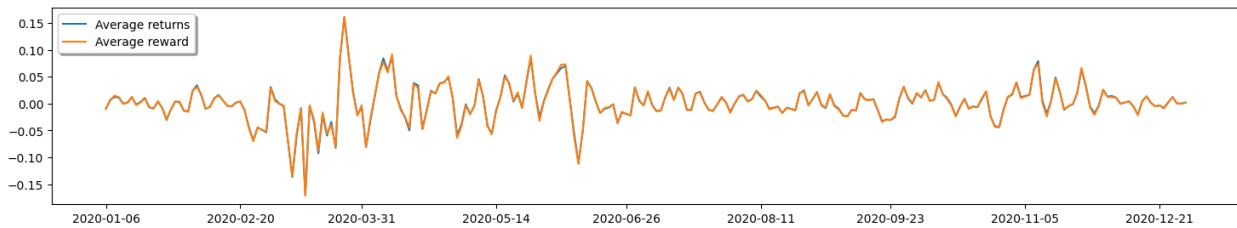


Figure 16: Average reward vs Average returns (UCB method)

```
Average frequency of optimal action = 0.0369
Holding the equally-weighted portfolio, average annualized return = 0.1515; average annualized standard deviation = 0.3865.
Holding the Bandit portfolio, average annualized return = 0.125; average annualized standard deviation = 0.3871.
```

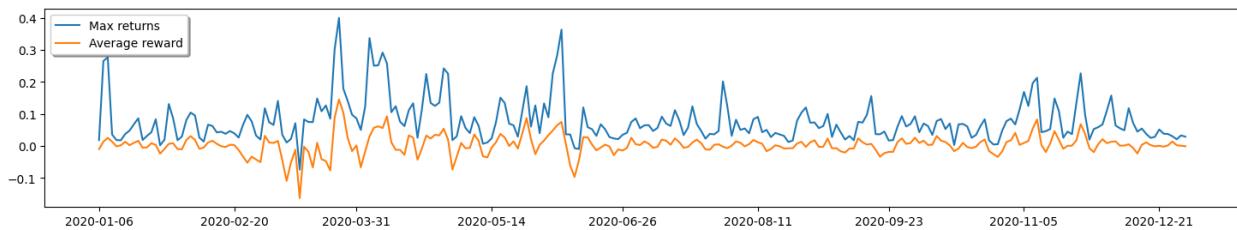


Figure 17: Average reward vs Max returns (Epsilon-greedy method)

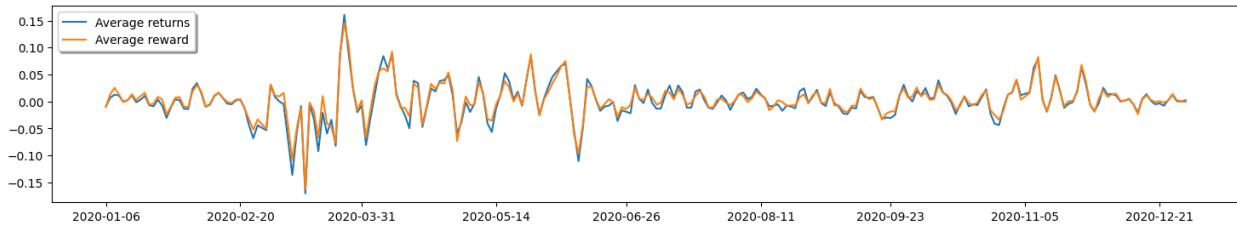


Figure 18: Average reward vs Average returns (Epsilon-greedy method)

```
Average frequency of optimal action = 0.0407
Holding the equally-weighted portfolio, average annualized return = 0.1515; average annualized standard deviation = 0.3865.
Holding the Bandit portfolio, average annualized return = 0.5006; average annualized standard deviation = 0.3489.
```

The most significant difference is shown in the simulation of investment initiating with \$1 in the portfolio. In the analysis using 2008 data, the ‘bandit’ strategy using both methods generates similar returns to the buy-and-hold strategy of the equally weighted portfolio. Nevertheless, when we rerun our analysis using more recent data (2020 data), the cumulative returns show a larger deviation. The ‘bandit’ strategy using the UCB method performs slightly poorer than the buy-and-hold strategy of the equally weighted portfolio while that using Epsilon-greedy method performs much better than the buy-and-hold strategy of the equally weighted portfolio.

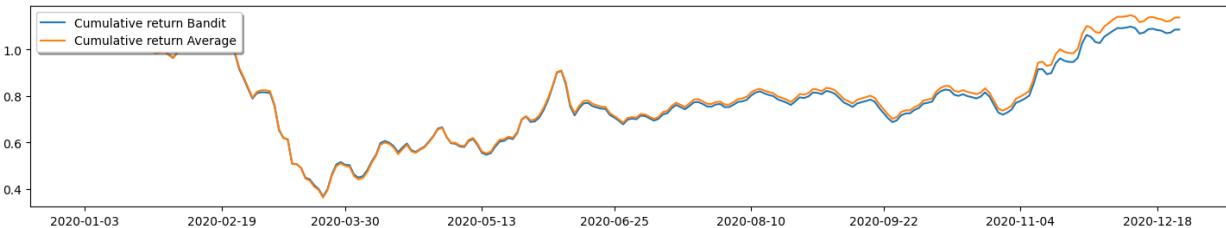


Figure 19: Cumulative return of Bandit portfolio vs equally-weighted portfolio (UCB method)

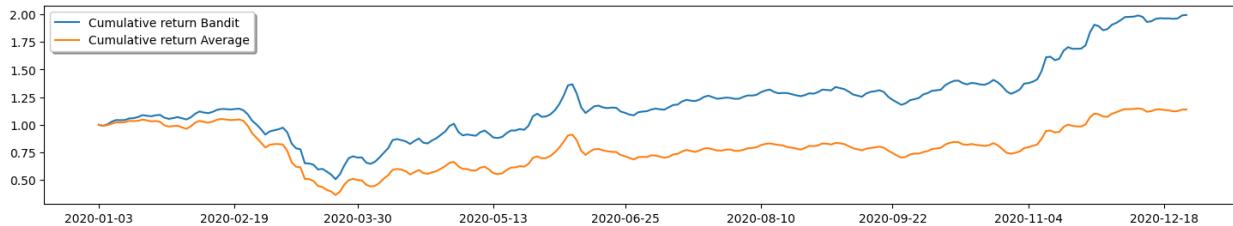


Figure 20: Cumulative return of Bandit portfolio vs equally-weighted portfolio (Epsilon-greedy method)

Although the use of recent data may indeed provide valuable insights and potentially enhance the performance of the algorithms employed, we cannot solely attribute the observed differences in results to the utilization of more recent data. While rerunning the algorithm, it is essential to acknowledge that other factors may have contributed to the variations in outcomes. For instance, modifying certain parameters as well as altering the length of the evaluation period can introduce additional elements of change.

References

1. Stochastic Modeling. WorldQuant University, 2023.
2. Huo. "Risk-aware multi-armed bandit problem with application to portfolio selection." 15 Nov. 2017, <https://royalsocietypublishing.org/doi/full/10.1098/rsos.171377>