

# Correlation Prediction

December 2023

## 1 Abstract

Our objective is to predict the correlation of GOOG, AMZN, JPM, GME, and XOM with the SPY index for the next two months using only historical prices. So we choose three methods to achieve this:

- Multivariate GARCH
- Multi-factors Model
- ARIMA-LSTM Hybrid

We conducted our analysis in a segmented manner, focusing on individual components separately. The complexity of the model we employed necessitated an extended period for parameter tuning. However, due to time constraints, we were unable to exhaustively fine-tune these parameters. Consequently, our results primarily reflect our initial approach and conceptual understanding. Given additional time, we believe we could undertake a more meticulous and detailed refinement of our model, potentially leading to enhanced outcomes.

## 2 Main Work

### 2.1 Multivariate GARCH

Multivariate GARCH models, namely models for dynamic conditional correlation (DCC), are what we need in this case. For a bi-variate random variable,

$$U_{t|t-1} = \begin{pmatrix} u_{1,t|t-1} \\ u_{2,t|t-1} \end{pmatrix}$$

with covariance matrix

$$\Sigma_{t|t-1} = \begin{bmatrix} \sigma_{1,t|t-1}^2 & \sigma_{12,t|t-1} \\ \sigma_{12,t|t-1} & \sigma_{2,t|t-1}^2 \end{bmatrix}$$

In addition, we define

$$\Psi_{t|t-1} = U_{t|t-1} U_{t|t-1}^T = \begin{bmatrix} u_{1,t|t-1}^2 & u_{1,t|t-1} u_{2,t|t-1} \\ u_{1,t|t-1} u_{2,t|t-1} & u_{2,t|t-1}^2 \end{bmatrix}$$

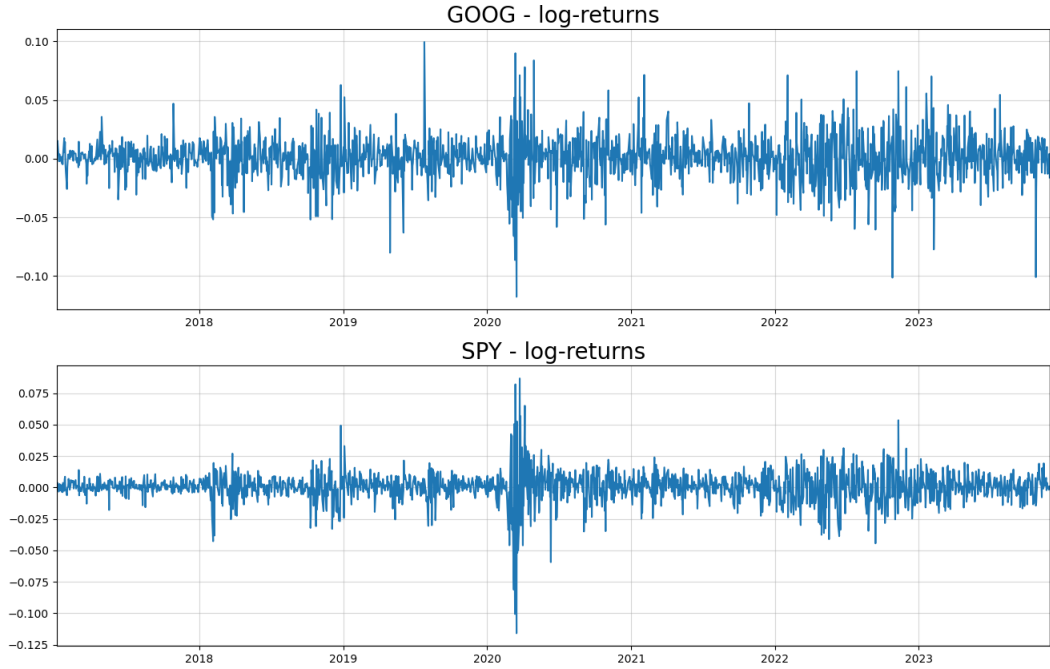
We condition covariance linearly on past covariances and past realizations of the actual random variables:

$$\Sigma_{t|t-1} = A_0 + A^T \Sigma_{t-1|t-2} A + B^T \Psi_{t-1|t-2} B$$

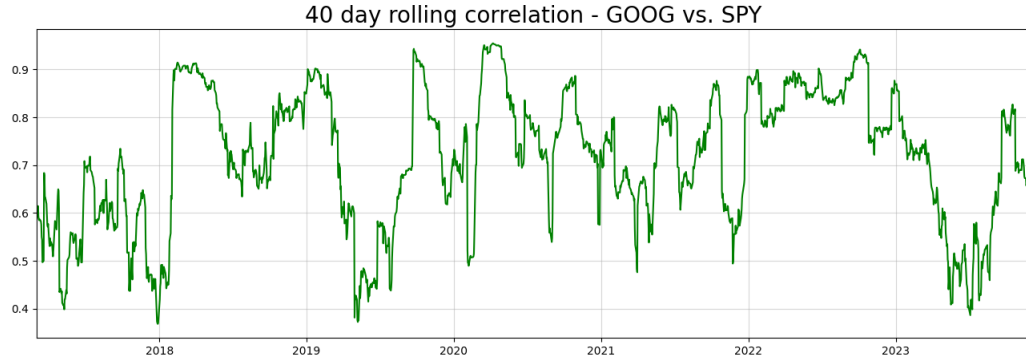
1. Calculate the in-sample distribution (get conditional dists(...)) is needed for optimization via maximum likelihood. This function calculates the likelihood values of each observation given the MGARCH model.

2. Forecast the out-of sample distribution (sample forecast(...)) - as the formulas for the model as a whole are quite complex, it's difficult to calculate the forecast distributions in closed form. However, we can, more or less, easily sample from the target distribution. With a sufficiently large sample, we can estimate all relevant quantities of interest (e.g. forecast mean and quantiles).

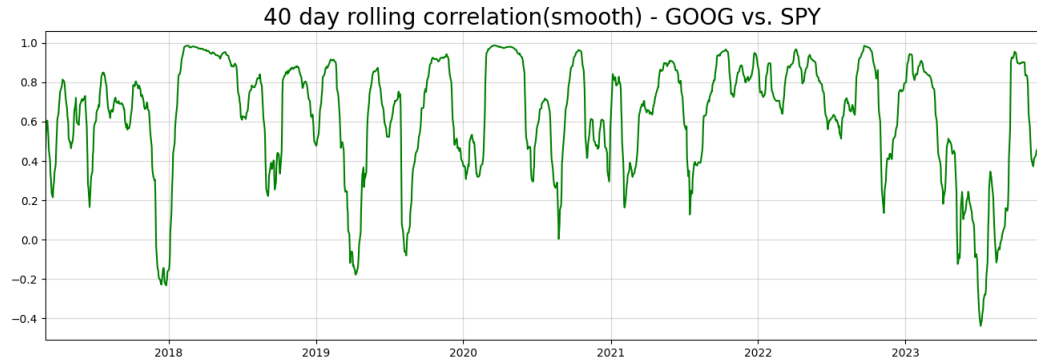
The log return of these two asset are shown below:



If we use rolling 40 days period (2 months) to see the correlation,



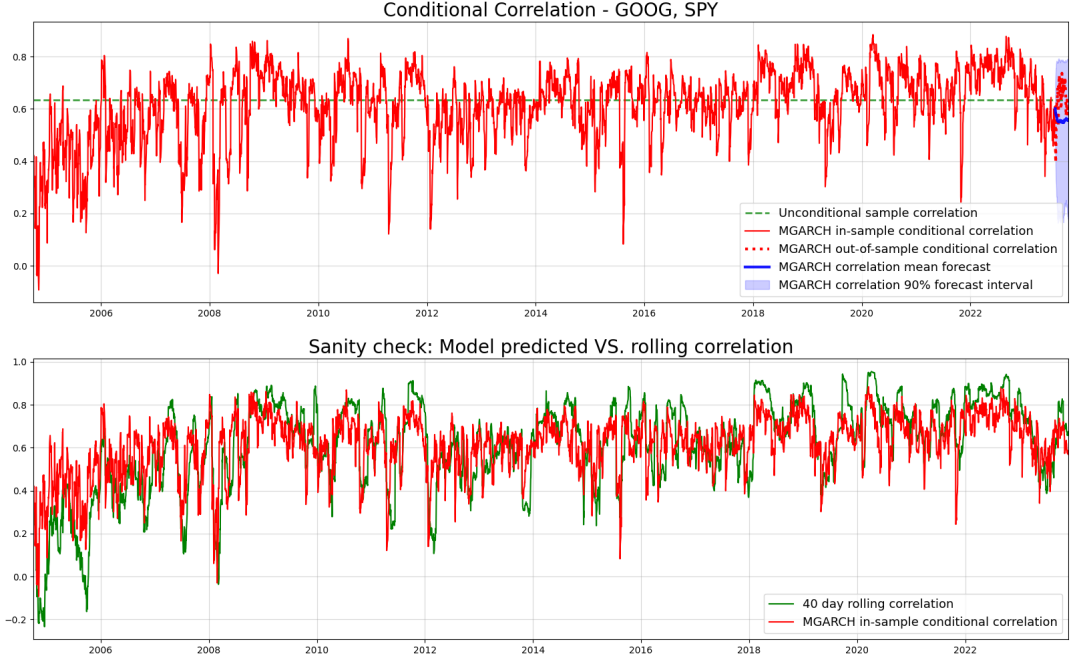
When we calculate the correlation for a 40-day rolling period, we find that the correlation curve changes abruptly at certain points (for example, from 0.6 one day to 0.9 the next). The predicting label is not smooth can go against our model training. So if we add an exponential Weighted Moving Average (EWMA) : EWMA gives a higher weight to recent observations and can capture recent market changes while also smoothing out excessive volatility.



It appears as if correlation between both indices has dropped since the beginning of the pandemic. Afterwards, correlation seems to fluctuate in cycles.

All in all, the pattern looks like a discretized version of an Ornstein-Uhlenbeck process. The error correction formulation in our model should be able to capture this behaviour accordingly.

We split the data into train and test set (last 90 observations), we can fit the model. Then we take samples from the (90 days ahead) forecast distribution as follows (this takes some time)



The forecasted correlation (blue) captures the actual correlation (red) under our model quite well. Obviously though, the true correlation is unknown. Nevertheless, our model matches the rolling correlation quite well, even out-of-sample. This implies that our approach is - at least - not completely off.

Also we plot the Sanity Check, to compare the rolling correlation and the conditional correlation that we modeled, we can see that they have the same trend at most of the times.

## 2.2 ARIMA-LSTM Hybrid Model

We apply LSTM recurrent neural networks (RNN) in predicting the stock price correlation coefficient of stocks and index. RNNs are competent in understanding temporal dependencies. The use of LSTM cells further enhances its long term predictive properties. To encompass both linearity and nonlinearity in the model, we adopt the ARIMA model as well. The ARIMA model filters linear tendencies in the data and passes on the residual value to the LSTM model. The ARIMA LSTM hybrid model is tested against other traditional predictive financial models such as the full historical model, constant correlation model, single index model and the multi group model. In our empirical study, the predictive ability of the ARIMA-LSTM model turned out superior to all other financial models by a significant scale. Our work implies that it is worth considering the ARIMA LSTM model to forecast correlation coefficient.

We still use rolling training method, first use ARIMA model to model the correlation between SPY and GOOG. Then put the residuals into LSTM model.

The model structure is shown below. The MAE between the predicted residual and actual residual is within 0.001.

But training LSTM for rolling takes so long time, so the result is omitted here.

Model: "sequential"		
Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 60, 100)	40800
dropout (Dropout)	(None, 60, 100)	0
lstm_1 (LSTM)	(None, 100)	80400
dropout_1 (Dropout)	(None, 100)	0
dense (Dense)	(None, 1)	101
Total params: 121301 (473.83 KB)		
Trainable params: 121301 (473.83 KB)		
Non-trainable params: 0 (0.00 Byte)		

## 2.3 Multi-factors Model

1. Data:

$$T_{train} + T_{valid} : 2000.1 - 2023.1$$

$$T_{test} : 2023.1 - 2023.11$$

2. Features:

Assuming that the set of technical factors is  $F$ , we use the technical factors today. These features contain information for the previous  $\tau$  days, where  $0 \leq \tau \leq 90$ .

3. labels:

Today is  $t$ , we calculate the correlation between stock and the index in next 40 days.

$$Label = corr(R_{stock:t:t+40}, R_{SPY:t:t+40})$$

4. Training: We use the expanding window method, In each iteration, the start point of the training window is fixed, but the end point moves forward, so the training set gradually grows larger.

5. Model Evaluation

Use MSE of correlation on  $T_{test}$  to evaluate the model.

### 2.3.1 Feature Construct

Firstly, we construct some technical signals to predict the next two months correlation. Take GOOG as example, We construct signals use both GOOG and SPY.

Table A.1: A description of all technical analysis features used

Name	Description	Formula
Simple Moving Average (SMA)	Simple moving average of the last $n$ observations of a time series $P^C$	$\frac{1}{n} \sum_{i=t-n}^t P_i^C, n \in \{6, 10, \dots, 50\}$
Exponential Moving Average (EMA)	Exponential moving average of a time series $P^C$	$\text{EMA}_t = P\alpha + \text{EMA}_t - 1(1 - \alpha)$ , where $\alpha = \frac{2}{1+N}, N \in \{6, 10, \dots, 50\}$
Bollinger Bands	Using the moving average, the upper and lower Bollinger bands are calculated as above and below 2 standard deviations of the closing price.	$\text{MA}_n \pm 2\sigma$ , $n \in \{14, 20, \dots, 60\}$
Momentum	Price change in the last $n$ periods	$P_t^C - P_{t-n}^C$ , $n \in \{6, 12, \dots, 60\}$
Acceleration	Difference in price change	$\text{Momentum}_t(n) - \text{Momentum}_{t-1}(n)$ , $n \in \{6, 12, \dots, 60\}$
Rate of Change	Rate of change of $P_t^C$	$\frac{P_t^C - P_{t-n}^C}{P_{t-n}^C} \cdot 100$ , $n \in \{4, 6, \dots, 60\}$
Moving Average Convergence Divergence (MACD)	Difference between two moving averages of slow and fast periods	$\text{EMA}_t(P^c, s) - \text{EMA}_t(P^c, f)$ , $s \in \{18, 24, 30\}$ , $f = 12$

Name	Description	Formula
Relative Strength Index	Compares the days that stock prices finished up against periods that stock prices finished down.	$100 - \frac{100}{1 + \frac{SMA_t(P_t^{up}, n_1)}{SMA_t(P_t^{dn}, n_1)}}$ , where $P_t^{up} = P_t^c$ if $P_t^c > P_{t-1}^c, P_t^{dn} = P_t^c$ if $P_t^c < P_{t-1}^c$ , 0 otherwise, and $n_1 \in \{8, 14, 20\}$
Stochastic Oscillator	Compares close price to a price range in a given period to establish if market is moving to higher or lower levels	$\frac{P_t^c - \min(P_n^{low})}{\max(P_n^{high}) - \min(P_n^{low})}$ , where $n \in \{10, 14, \dots, 22\}$
Williams Indicator	Captures moments when the market is overbought or oversold by calculating index	$\frac{\max(P_n^{high}) - P_t^c}{\max(P_n^{high}) - \min(P_n^{low})}$ , where $n = 14$
Money Flow Index	Measures the strength of money flow in and out of a stock.	$100 - \frac{100}{1 + \frac{PMF_t(n)}{NMF_t(n)}}$ , where $n = 14$ , $MF_t = P_t^{yp} \cdot VOL_t$ , and $PMF_t(n) = SMA_t(MF_t, n)$ when $MF_t > 0$ , $NMF_t(n) = SMA_t(MF_t, n)$ when $MF_t < 0$
Chaikin volatility	Evaluates the widening of the range between high and low prices.	$\frac{EMA(P^h - P^l, n)}{EMA_{t-n_1}(P^h - P^l, n)} - 1$ , where $n = n_1 = 10$
Negative and Positive volume index	Signals of bull markets based on the volume traded and race of change of prices.	$NVI_t = NVI_{t-1} + ROC_t(n)NVI_{t-1}$ if $VOL_t < VOL_{t-1}$ , and $NVI_{t-1}$ otherwise. $PVI_t = PVI_{t-1} + ROC_t(n)PVI_{t-1}$ if $VOL_t > VOL_{t-1}$ , and $PVI_{t-1}$ otherwise.
Price Volume Trend	Calculates cumulative total of volume where portion of volume added/subtracted is given by increase of decrease of close price with respect to previous period	$\sum_{t=1}^n VOL_t \cdot ROC_t(n_1)$ , where $n_1 = 1$
On Balance Volume	Evaluates impact of positive and negative volume flows	$OBV_t = OBV_{t-1} \pm VOL_t$ when $P_t^c$ is greater than or less than $P_{t-1}^c$ respectively
Accumulation/Distribution line	Evaluates the effect of accumulative flow of money. Significant differences between the accumulation distribution line and the price produce signals.	$\sum_{t=1}^n CLV_t \cdot VOL_t$ , where $CLV_t = \frac{2P_t^{uc} - P_t^l - P_t^h}{P_t^h - P_t^l}$

Also, beside these we also add some features to describe the correlation and interaction between the stock and the index.

Such as different period of correlation, and the technical signals to analyze the correlation.

Then we merge the signals and fill the NA values and transform them by Min-max scaler.

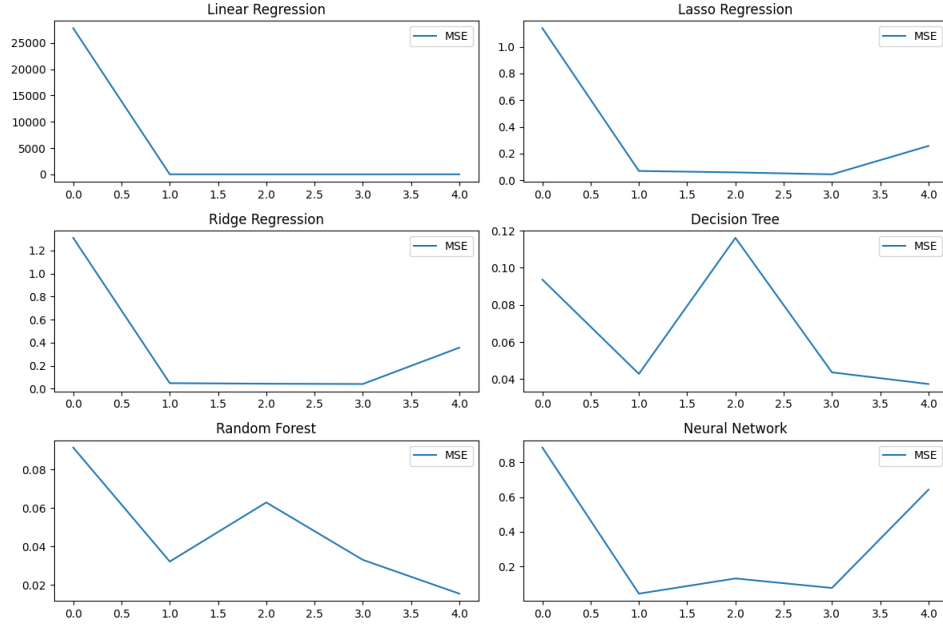
### 2.3.2 Training Procedure

We use time-series cross validation method. TimeSeriesSplit is a function provided by the scikit-learn library in Python, specifically designed for time series cross-validation. This function generates train/test indices to split time series

data samples in a sequential manner. It is part of the model selection module of scikit-learn.

We choose the data set before 2023-01-01 as training and validation set(keep expanding the training set and take the following as validation set.)

We split the training and validation set into 5 folds and see the MSE plots in five training sets.



The best model is Random Forest, and best model MSE on Testing Set is 0.04508583990479208. On average, there is a difference of 0.2 between our predicted value and the actual value. The predicted value in 2023 is shown below.

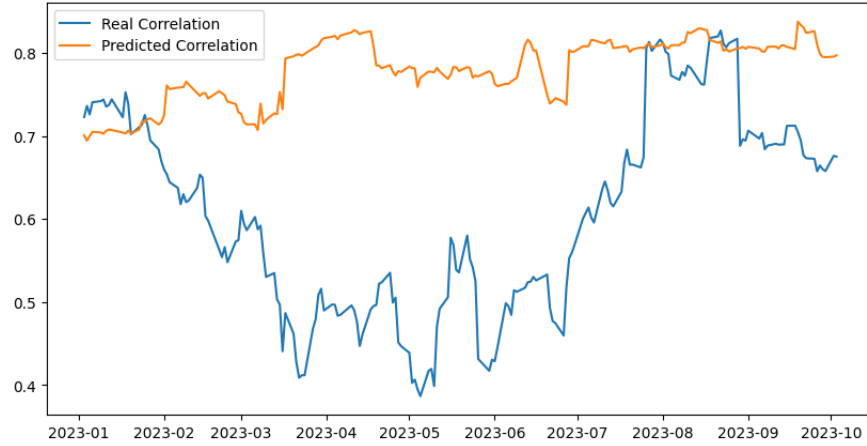
Model	validation MSE	R2
Linear Regression	5559.493834	0.913144
Lasso Regression	0.314127	0.663563
Ridge Regression	0.360557	0.653635
Decision Tree	0.072743	0.999246
Random Forest	0.045931	0.996539
Neural Network	0.154257	0.784062

Table 1: Mean Squared Error (MSE) and R2 for different models

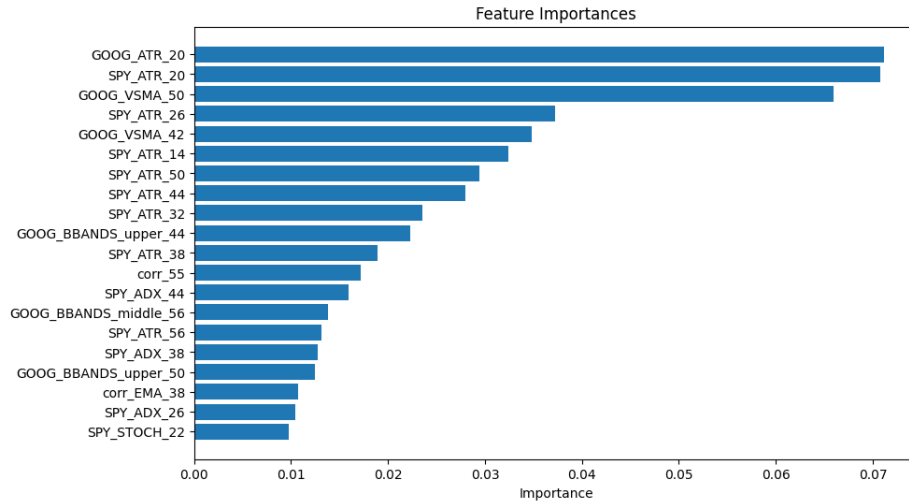
We can see that the linear regression on the small training set has large MSE on validation set. Considering that many of our features have high correlation,



we choose the regularization method such as lasso and ridge, to solve the multicollinearity.



Also we check the feature importance of the features and keep the first 15 features:



The Average True Range (ATR) is a technical analysis indicator that measures market volatility by decomposing the entire range of an asset price for that period. The ATR is calculated using the true range over a specified period, typically 14 days.

Components of ATR:

1. True Range (TR): The True Range for a period is the greatest of the following:

- Current High less the current Low
  - The absolute value of the current High less the previous Close
  - The absolute value of the current Low less the previous Close
2. Average True Range (ATR): ATR is the moving average of the TR over the specified period. It can be calculated using various moving average types, but the most common is the simple moving average.

Calculation: Given a DataFrame 'df' with columns 'High', 'Low', and 'Adj Close', and a time period 'i', the ATR is calculated as follows:

1. Calculate the True Range (TR) for each period:

$$TR = \max[(High - Low), \text{abs}(High - PreviousClose), \text{abs}(Low - PreviousClose)]$$

2. Calculate the Average True Range (ATR):

ATR = Moving Average of TR over the last 'i' periods.

The VSMA is calculated by taking the sum of the product of the volume and price for each period, divided by the sum of the volume for that period.

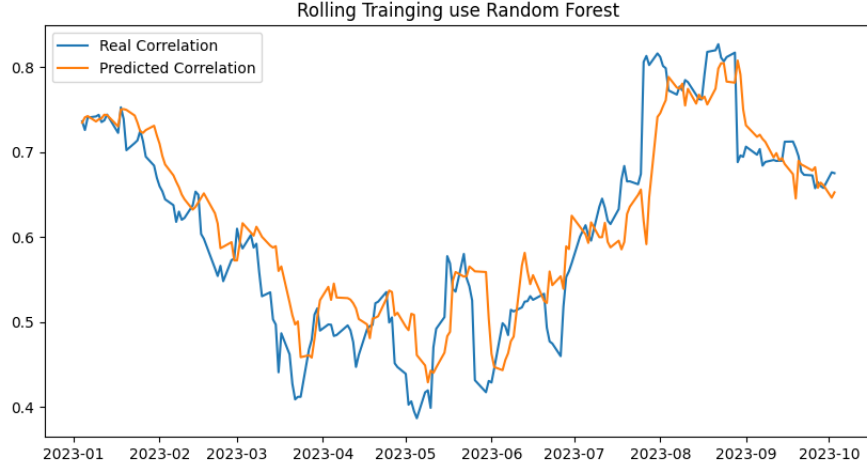
$$VSMA_t = \frac{\sum_{i=t-n+1}^t (Price_i \times Volume_i)}{\sum_{i=t-n+1}^t Volume_i}$$

In this formula:

- $VSMA_t$  is the Volume Weighted Moving Average at time  $t$ .
- $Price_i$  is the price at time  $i$ .
- $Volume_i$  is the volume at time  $i$ .
- $n$  is the number of periods over which the average is computed (the window size).
- The summation runs from  $t - n + 1$  to  $t$ , summing over the last  $n$  periods.

### 2.3.3 Rolling Training

In our comprehensive analysis, we implemented a rolling training methodology to predict asset correlations, utilizing a Random Forest model. This approach was predicated on the hypothesis that more recent data would yield more accurate predictions due to the dynamic nature of financial markets. We systematically varied the rolling window periods to optimize our model's performance. Intriguingly, we observed a consistent improvement in prediction accuracy with the adoption of rolling training as opposed to static historical data training. After extensive parameter tuning, a rolling period of 60 days emerged as the most effective, striking a balance between capturing recent market trends and maintaining sufficient data for robust model training. This was substantiated by a comparative analysis of the Mean Squared Error (MSE) across various window sizes, where the 60-day window consistently outperformed others, indicating its superior predictive capability in our modeling context.

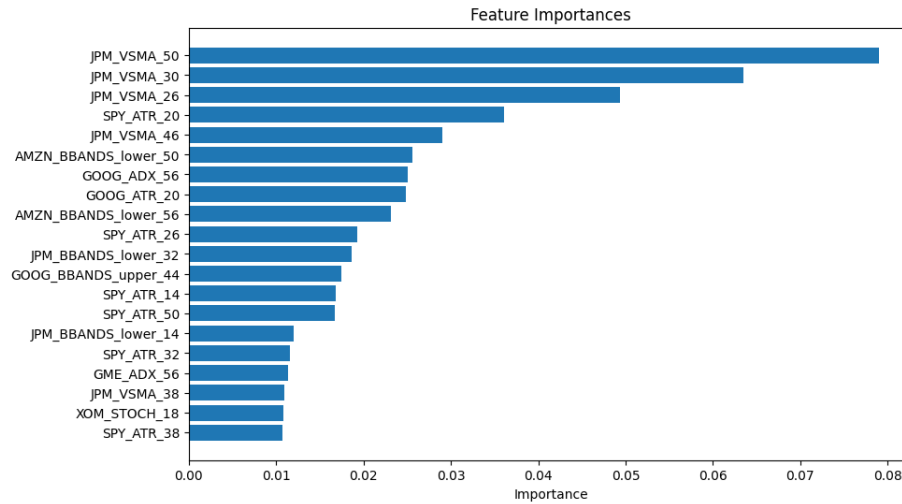


Rolling Period (Days)	MSE
30	0.005835
40	0.003145
50	0.002934
60	0.002626
70	0.002677
80	0.002645

Table 2: Mean Squared Error (MSE) in validation set for Different Rolling Periods

#### 2.3.4 Features from other stocks

We are considering that does the signals from other stocks also impacted the current stock? So we add features from other stocks. And retrain the Random Forest Model and show that the MSE in validation set decreasing about 0.01.(Average MSE: 0.034). And the feature importance rank also changed alot. We can see the importance rank that:



The features constructed from JPM contributes a lot to the prediction of correlation with GOOG and SPY!

The insightful analysis reveals that the features originating from JPM play a pivotal role in accurately forecasting the correlation dynamics between GOOG and the SPY. This finding underscores a noteworthy market interplay, indicating that the correlation existing between an individual stock and a broad market index can be substantially affected by the movements or characteristics of another stock. This highlights the complex and interconnected nature of financial markets, where the performance and characteristics of one key player, such as JPM, can exert a measurable influence on the correlation patterns between other major entities like individual stocks and market indices.

The other stock's results will be contained in code.