

## Forecasting Realized Volatility of the Oil Future Prices via Machine Learning<sup>\*</sup>

Byung-June Kim<sup>a</sup>, Taeyoon Kim<sup>b</sup>, Myung-Jun Kim<sup>c, †</sup>, and Bong-Gyu Jang<sup>d, †</sup>

<sup>a</sup>Department of Industrial and Management Engineering, POSTECH, Korea. Tel: +82-54-279-2864, Fax: +82-54-279-2870, [kby219@postech.ac.kr](mailto:kby219@postech.ac.kr)

<sup>b</sup>Department of Industrial and Management Engineering, POSTECH, Korea. Tel: +82-54-279-2864, Fax: +82-54-279-2870, [taeyoonkim.ime2016@postech.ac.kr](mailto:taeyoonkim.ime2016@postech.ac.kr)

<sup>c, †</sup>Corresponding Author, Department of Industrial and Management Engineering, POSTECH, Korea. Tel: +82-54-279-2864, Fax: +82-54-279-2870, [kmj4720@postech.ac.kr](mailto:kmj4720@postech.ac.kr)

<sup>d, †</sup>Corresponding Author, Department of Industrial and Management Engineering, POSTECH, Korea. Tel: +82-54-279-2864, Fax: +82-54-279-2870, [bonggyujang@postech.ac.kr](mailto:bonggyujang@postech.ac.kr)

### ABSTRACT

This paper explores the potential use of machine learning models in crude oil realized volatility forecasting through a variety of empirical analyses and robustness checks. Although the conventional Heterogeneous Autoregressive (HAR) model is widely accepted, the machine learning models with the HAR factors can significantly improve the forecasting performance. We also found that macroeconomic variables such as supply factors, implied volatility indices and uncertainty factors can be useful in forecasting oil volatility.

### KEYWORDS

Volatility Forecasting; Oil Future; Machine Learning; Forecasting Model

### JEL CLASSIFICATION

C5, C22, G1, Q4

---

<sup>\*</sup>This work was supported by the Ministry of Education of the Republic of Korea and the National Research Foundation of Korea.(NRF-2022S1A3A2A02089950)

## 1. Introduction

Volatility is a key characteristic of commodity markets, particularly evident in the oil sector. Historically, oil prices have been subject to considerable fluctuations, driven by a variety of factors including geopolitical unrest, imbalances in supply and demand, market speculation, and major economic changes (Hamilton, 2009). The advent of the COVID-19 pandemic has exerted a profound impact on the commodity market, leading to "unprecedented disruptions in worldwide supply chains, dramatic shifts in oil demand, and increased economic instability."

These events have triggered severe swings in oil prices, influencing a wide spectrum of market stakeholders from traders and investors to policymakers. The unparalleled drop in demand, severe conditions in the job market, and restrained consumer spending resulting from the COVID-19 crisis are unparalleled in living history. A stark example of this was when dwindling demand drove oil future prices into negative figures, reaching -\$36.98 in May 2020.

The extraordinary plunge in oil prices is not to be underestimated due to the considerable economic repercussions it carries for industries reliant on oil (Apergis and Miller, 2009; and Kilian, 2009).<sup>1</sup> Even prior to the global ramifications of COVID-19, uncertainty around oil prices had been on a gradual rise over the past decades, posing a substantial risk to the global economy.

This paper explores the possibility of the potential usage of machine learning in the field of volatility forecasting. Comparison of various forecasting models with a rich set of data and various forecasting horizons guides us to the fact that a combination of conventional models and machine learning techniques can improve the forecasting performance in out-of-sample. Ma et al. (2017) and Ma et al. (2018) investigated autoregressive-type forecasting models for oil price volatility, including the heterogeneous autoregressive (HAR) models. However, they covered only autoregressive-type models and did not check the time consistency of the forecasting performance. [Unfortunately, although](#)

---

<sup>1</sup>The 1980s oil glut, for example, has caused a blow to the exports of the Soviet Union, and it may have been a factor in accelerating the collapse of the Soviet Union. On the other side, many crop prices also remained in the doldrums amid weak raw material prices, which resulted in a recession in the U.S. agricultural industry, leading to bankruptcy in many small banks. As another example, the increased U.S. shale supply caused a sharp fall in oil prices again in 2014, which has yet to recover since then. As a result, Russia fell into recession, and growth in the Middle East slowed.

HAR-X type models are successful in-sample prediction, the performance of HAR-X is not as effective as that of other machine learning models in out-of-sample. We show that the machine learning models with factors in the HAR model as input variables have the potential to enhance the out-of-sample forecasting performance with one-week, bi-week, and month ahead horizons.

We examine the out-of-sample forecasting performance of ten different models as well as their robustness with four different metrics:  $\hat{R}_{oos}^2$  score,  $R_{oos}^2$  score, Model Confidence Set (MCS) test (Hansen et al. (2011)), and Diebold-Mariano (DM) test (Diebold and Mariano, 2002). We consider the period from April 2002 to April 2024 to include the Great Recession and COVID-19.

The statistic  $\hat{R}_{oos}^2$  is employed to compare the forecasting residual with the forecasting residual of the benchmark model, designated as HAR. The application of machine learning models incorporating HAR factors, such as random forest regression (RFR) or ElasticNet, has the potential to achieve notably high values of  $R_{oos}^2$  in out-of-sample. We found that the longer forecast horizon tends to yield more accurate forecasts of models, which may be attributed to the presence of noise due to data frequency. Additionally, the performance orders can vary across different forecast horizons. The performance of conventional HAR and HAR-X models illustrates the distinct importance of lagged variables for different horizons. Exposure to varying levels of autocorrelation effects may contribute to the observed performance order differences across different horizons.

While  $\hat{R}_{oos}^2$  is based on the benchmark model,  $R_{oos}^2$  compares the forecasting residual with the variance of true values in out-of-sample (Guo and Lin (2020)). It measures the predictability of the model by residuals between predicted values and true values. With  $R_{oos}^2$ , we can compare the forecasting performance of different lengths of the test period. In our study, the performance of the conventional HAR model decreases as the period becomes more recent. This might suggest that the importance of factors other than lagged variables, that is, memory or momentum variables, has increased recently. In contrast, some machine learning models such as the random forest regression and the recurrent neural network model, show very stable  $R_{oos}^2$  scores even during the period of COVID-19.

In addition, we construct a rich set of data including implied volatility indices and

uncertainty factors, following Degiannakis and Filis (2022), Delis et al. (2022), Delis et al. (2023), Miao et al. (2017), Wei et al. (2017), and Ma et al. (2018).<sup>2</sup> As expected in Ma et al. (2018), appending uncertainty factors help models improve the out-of-sample performance and its robustness. Furthermore, the number of selected features plays a crucial role in performance. Linear models like HAR-X, LASSO, and ElasticNet tend to perform better as more features are included. In contrast, nonlinear machine learning models such as RFR and GBR achieved high out-of-sample scores with just 20 features, suggesting that too many features could negatively impact their performance. Deep learning models, including ANN and RNN, show more variable results. This might be due to the nonlinear relationship between realized volatility of crude oil and uncertainty indices.

Our study is distinct from other research in some significant respects. First, we incorporated uncertainty factors, which permitted us to examine a more expansive spectrum of variables that influence oil volatility. Secondly, we employed feature selection techniques with diverse machine learning models to examine the extensive range of models and the vast number of features. Moreover, our models are all interpretable, enabling us to ascertain which feature plays a pivotal role in predicting oil prices. Finally, our methods produced favorable outcomes, demonstrating the efficacy of our approach. These concepts offer invaluable insights and a more efficient research methodology than previous studies.

The remaining sections of the paper are organized as follows. In section 2, we illustrate input data and data preprocessing procedure. In section 3, we explain the training procedure of various forecasting models. In section 4, we present the out-of-sample  $\hat{R}_{oos}^2$  score,  $R_{oos}^2$  score, MCS test, and DM test, as the forecasting performance measures, and analyze the results of the forecasting models. We also present which feature is chosen in the models frequently.

Factor Group	Individual Variable	Frequency	Source
Prices	WTI (West Texas Intermediate) future prices	Daily	Energy Information Administration
	WTI spot prices	Daily	Energy Information Administration
	Brent oil spot prices	Daily	Energy Information Administration
	NGL (Natural Gas Liquids) future prices	Daily	Energy Information Administration
	NGL spot prices	Daily	Energy Information Administration
Supply Factors	Global crude oil production	Monthly	JODI-Oil Database
	Global crude oil stock	Monthly	JODI-Oil Database
	Global crude oil export	Monthly	JODI-Oil Database
	Total OPEC production capacity	Monthly	Energy Information Administration
	Capacity utilization rate	Weekly	Energy Information Administration
Demand Factors	Global crude oil import	Monthly	JODI-Oil Database
	Liquid fuels consumption in World	Monthly	Energy Information Administration
	PPI in China	Monthly	Federal Reserve Bank of St. Louis Economic Database
	PPI in US	Monthly	Federal Reserve Bank of St. Louis Economic Database
	PPI in EU	Monthly	Federal Reserve Bank of St. Louis Economic Database
Financial Factors	S&P 500 Adjusted Close	Daily	Yahoo Finance
	Japan / US Foreign Exchange Rate	Daily	Federal Reserve Bank of St. Louis Economic Database
	US / Euro Foreign Exchange Rate	Daily	Federal Reserve Bank of St. Louis Economic Database
	US / UK Foreign Exchange Rate	Daily	Federal Reserve Bank of St. Louis Economic Database
	China / US Foreign Exchange Rate	Daily	Federal Reserve Bank of St. Louis Economic Database
	Federal Funds Rate	Monthly	Federal Reserve Bank of St. Louis Economic Database
	MSCI World Standard (Large+Mid Cap)	Monthly	MSCI
Implied Volatility Indices	CBOE Volatility Index	Daily	Yahoo Finance
	CBOE Crude Oil Volatility Index	Daily	Yahoo Finance
	CBOE DJIA Volatility Index	Daily	Yahoo Finance
	CBOE Gold Volatility Index	Daily	Yahoo Finance
Uncertainty Factors	Global Economic Policy Uncertainty (current)	Monthly	Davis (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	Global Economic Policy Uncertainty (ppp)	Monthly	Davis (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	Daily Infectious Disease Equity Market Volatility Tracker	Daily	Baker et al. (2019) and Baker et al. (2020) <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Economic Policy Uncertainty	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Monetary Policy	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Fiscal Policy (Taxes or spending)	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Taxes	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Government Spending	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Healthcare	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in National Security	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Entitlement Programs	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Regulation	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Financial Regulation	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Trade Policy	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	US Economic Policy Uncertainty in Sovereign Debt and currency crises	Monthly	Baker et al. (2016), <a href="https://www.policyuncertainty.com">https://www.policyuncertainty.com</a>
	World Uncertainty Index	Quarterly	Ahir et al. (2018), <a href="https://worlduncertaintyindex.com/">https://worlduncertaintyindex.com/</a>
	World Pandemic Uncertainty Index	Quarterly	Ahir et al. (2018), <a href="https://worlduncertaintyindex.com/">https://worlduncertaintyindex.com/</a>

**Table 1.** There are six groups of data - prices, supply factors, demand factors, financial factors, implied volatility indices, and uncertainty factors. The dependent variable is WTI future price, and the rest of the data are explanatory variables used in the models.

## 2. Data

In line with Miao et al. (2017), Wei et al. (2017), and Ma et al. (2018), we construct the realized volatility of oil future contracts on the West Texas Intermediate (WTI) crude oil as follows:

$$RV_t = \sqrt{\sum_{j=1}^{N_t} r_{t,j}^2}, \quad \text{for } r_{t,j} = 100 \times \log(p_{t,j}/p_{t,j-1}),$$

where  $N_t$  indicates the number of business days in the  $t$ -th week, and  $p_{t,j}$  represents the daily WTI future prices on  $j$ -th business day of the  $t$ -th week. We consider the period from April 2002 to April 2024 to include the Great Recession and COVID-19. Because each future contract must be exercised at maturity, to estimate continuous volatility, daily prices are calculated based on the roll-over rule (Schwager, 2017).<sup>3</sup> The set of

<sup>2</sup>Kilian and Hicks (2013) argued that repeated shocks of oil demand cause the oil price shock. Other plentiful researches also assert interdependence between the oil price uncertainty and the structural fluctuation of global economic activity (Herrera et al., 2019; and Caggiano et al., 2020).

<sup>3</sup>We follow ‘Definitions, Sources, and Explanatory Notes’ in the EIA webpage ([https://www.eia.gov/dnav/pet/TblDefs/pet\\_pri\\_fut\\_tbldef2.asp](https://www.eia.gov/dnav/pet/TblDefs/pet_pri_fut_tbldef2.asp)) to calculate continuous WTI future prices. On each monthly roll-over

uncertainty indices is more closely associated with macroeconomic trends than with the specific behavior of the oil market. Consequently, a sudden change in these indices will require a period of adjustment in the oil market. We hypothesize that uncertainty indices do not exert an immediate (daily) influence on oil prices, and thus, we anticipate weekly price movements.

A total of 42 explanatory variables were employed in the analysis (see Table 1). Each explanatory variable is classified into one of six categories: prices, supply factors, demand factors, financial factors, implied volatility indices, and uncertainty factors. We apply a difference transformation to the data that does not satisfy stationarity criteria over the entire period. For higher-order differences (greater or equal to 3), we set the order of difference to 2 to preserve information. Additionally, we consider the publication lags of the explanatory variables to mitigate potential look-ahead bias in the forecast.

**Prices** West Texas Intermediate (WTI) has a long-run equilibrium relationship with Brent (Hammoudeh et al., 2008), and close interconnection with natural gas (Brown and Yucel, 2008). We adopted the future and spot prices of WTI (only spot prices), Brent, and natural gas liquid (NGL) from the US Energy Information Administration (EIA), as explanatory variables.

**Supply** Global crude oil data, including production, stock, and export, are from world-primary data in the JODI-Oil Database. Oil supply and demand are accepted factors to describe the oil price dynamics (Kilian, 2009; Ma et al., 2018; Wei et al., 2017). Hallock Jr et al. (2004) showed a strong relationship between oil production and exportation. Hamilton (2009) asserted the importance of oil inventory to prices. Total OPEC production capacity and Capacity utilization rate are from EIA.

**Demand** Demand factors have a significant influence on oil prices (Hamilton, 2009; and Kilian, 2009) The demand factors consist of the global crude oil import from the JODI-Oil Database, liquid fuel consumption in the World from the EIA, and the producer Price Index (PPI) in China, the US, and the EU from the Federal Reserve Bank (FRB) of St. Louis Economic Database.

---

day, the oil price is adjusted by the price difference between the future contract with the first-nearest-to-maturity and the future contract with the second-nearest-to-maturity. The weekly realized volatility is calculated by the squared sum of daily WTI returns on each last business day of the week. Each contract expires on the third business day before the 25th calendar day. When the 25th calendar day is not a business day, the contract expires on the third business day before the latest business day prior to the 25th calendar day. We excluded the date when the oil price became negative. More details about the roll-over rule are in Appendix A.

**Financial Factors** Following Miao et al. (2017), financial factors include the federal fund rate, the daily exchange rate of Japan-US, US-EU, US-UK, and China-US from the FRB, the S&P 500 Adjusted Close data from Yahoo Finance and MSCI World Standard (Large+Mid Cap) from MSCI.

**Implied Volatility Indices** As showed in the Delis et al. (2022) and Delis et al. (2023), the implied volatility indices have potential explanatory power in forecasting the oil implied volatility index. In this regard, four implied volatility indices are included: the CBOE Volatility Index (VIX), the CBOE Crude Oil Volatility Index (OVX), the CBOE DJIA Volatility Index (VXD), and the CBOE Gold Volatility Index (GVZ). The data were collected from Yahoo Finance.

**Uncertainty Factors** Ma et al. (2018) also considers uncertainty indices when handling the oil. Following Ma et al. (2018), Baker et al. (2016), Baker et al. (2019), Baker et al. (2020), Ahir et al. (2018), and Davis (2016) with their websites, we use various uncertainty indices as an input.<sup>4</sup>

### 3. Forecasting Models

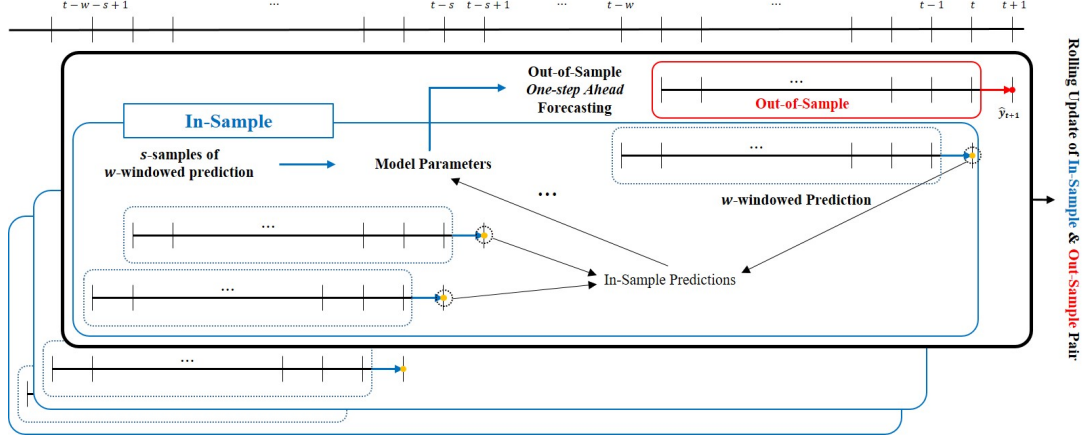
#### 3.1. Cross Validation

In the time-series model, preventing the use of future information is a pertinent first step. Similar to Gu et al. (2020), we construct the out-of-sample data and design cross-validation as shown in Fig 1. For each time, we stack  $s$  historical data, each of which has length- $w$  training data. With  $s \times w$  training input and corresponding  $s \times 1$  output, standard cross-validation is performed to find the best hyperparameters.

Now one-step ahead length- $w$  data (test data or the out-of-sample data) is chosen to be a test input. With the in-sample trained model and the one-step-ahead length- $w$  input, the model predicts one-week volatility. After prediction, we reset the model and move our  $s$  set of historical data forward, and follow the same procedure. This process guarantees that models do not use future information as input. The number of samples

---

<sup>4</sup>Economic uncertainty data were retrieved from <https://www.policyuncertainty.com/> and <https://worlduncertaintyindex.com/>. This webpage includes data from numerous papers covering economic policy uncertainty data (EPU), categorical economic policy uncertainty data, world uncertainty index (WUI) to even infectious disease data.



**Figure 1. Walk-Forward Cross-Validation** A single forecast on dependent variable  $y_{t+1}$  in time  $t + 1$  is based on explanatory variables in the time interval  $[t - w + 1, t]$ . Regardless of the forecasting model choice, this single forecast needs parameters for each explanatory variable in the  $w$ -windowed time interval. The estimation of the corresponding parameters for each explanatory variable uses  $s$ -sample points by minimizing the objective function, e.g., the sum of the least square, of in-sample residuals. When  $w$  is large, a single forecast considers more past information. A large number of samples  $s$  would stabilize parameter estimation, but it would catch up with recent changes slowly.

( $s$ ) and windows ( $w$ ) is chosen to contain enough information from the past but not overfit to the training data (Fig 1).<sup>5</sup>

### 3.2. Forecasting Models

Machine learning models are becoming more prevalent in the field of forecasting (Gu et al., 2020; Ghoddusi et al., 2019). We use both time series models - the heterogeneous autoregressive with exogenous variables (HAR-X), the time-varying parameter heterogeneous autoregressive (TV-HAR) - and various machine learning models.<sup>6</sup> The HAR-X model from Corsi (2009) is a benchmark. The reader can find the analysis on various machine learning models - the least absolute shrinkage and selection operator (LASSO), the elastic net (EN), the decision tree regression (DTR), the random forest regression (RFR), the gradient boosting regression (GBR), the artificial neural network (ANN), and the recurrent neural network (RNN). Detailed descriptions of how we train models are in Appendix B.

<sup>5</sup>In this study, all models use  $w = 1, s = 52$ .

<sup>6</sup>The reader should beware that depending on the operating system, version of packages, type of CPUs, and other factors may affect the results. We have tested in different environments, and in most cases, the differences were relatively insignificant. Additionally, alternative models, including ARIMA-X and principal component regression (PCR), were evaluated. However, the outcomes were unsatisfactory, and thus, they have been excluded from this paper.



## Heterogeneous Autoregressive with Exogenous Variables

Corsi (2009) found the stylized facts about the long memory and multi-scaling behavior of the realized volatility. The HAR-X contains two explanatory variables: weekly and monthly realized volatility. These variables represent short-term and medium-term, respectively.<sup>7</sup> The HAR-X is defined by

$$\begin{aligned}
\widehat{RV}_{t+1}^W &= c_W + \beta_{W,w}RV_t^W + \beta_{W,m}RV_{t-4:t}^W + \beta_{W,q}RV_{t-13:t}^W \\
&\quad + \sum_{k=0}^{w-1} \sum_{m=1}^n \eta_{W,m,t-k} x_{m,t-k} + \epsilon_{W,t+1} \\
\widehat{RV}_{t+1}^{2W} &= c_{2W} + \beta_{2W,w}RV_t^{2W} + \beta_{2W,q}RV_{t-6:t}^{2W} + \beta_{2W,y}RV_{t-26:t}^{2W} \\
&\quad + \sum_{k=0}^{n-1} \sum_{m=1}^n \eta_{2W,m,t-k} x_{m,t-k} + \epsilon_{2W,t+1} \\
\widehat{RV}_{t+1}^M &= c_M + \beta_{M,w}RV_t^M + \beta_{M,q}RV_{t-3:t}^M + \beta_{M,y}RV_{t-12:t}^M \\
&\quad + \sum_{k=0}^n \eta_{m=1}^{m,t-k} x_{m,t-k} + \epsilon_{M,t+1}
\end{aligned}$$

where  $RV_t^W$ ,  $RV_t^{2W}$ , and  $RV_t^M$  are the weekly, bi-weekly, and monthly realized volatility based on daily return, respectively. For the weekly realized volatility  $RV_t^W$ , the monthly average realized volatility  $RV_{t-4:t}^W$  and the quarterly realized volatility  $RV_{t-13:t}^W$ , which are the average  $RV$  from  $t-4$  to  $t$  and  $t-13$  to  $t$  respectively, represent the short-term variable and long-term variable. Every model below also contains realized volatility factors as explanatory variables. The parameter estimation process is the same as the linear regression.<sup>8</sup>

## Time-Varying Parameter Heterogeneous Autoregressive

The time-varying parameter heterogeneous autoregressive (TV-HAR) is an extension of the HAR, considering the parameter changes over time. This method is commonly used to enhance the forecasting power of the HAR-type models (Delis et al. (2022)). In

---

<sup>7</sup>The original HAR contains daily realized volatility to forecast daily realized volatility. However, our goal in this paper is to forecast the weekly realized volatility. So, we exclude the daily realized volatility factor.

<sup>8</sup>Therefore, machine learning models are the extended version of the HAR, in the aspect of variable usage. When the machine learning models do not consider realized volatility factors as explanatory variables, they perform poorly.

TV-HAR, the coefficients in the equation above change to time-varying coefficients as follows:

$$\begin{aligned}
\widehat{RV}_{t+1}^W &= c_W + \beta_{W,w,t}RV_t^W + \beta_{W,m,t}RV_{t-4:t}^W + \beta_{W,q,t}RV_{t-13:t}^W \\
&\quad + \sum_{k=0}^{w-1} \sum_{m=1}^n \eta_{W,m,t-k}x_{m,t-k} + \epsilon_{W,t+1} \\
\widehat{RV}_{t+1}^{2W} &= c_{2W} + \beta_{2W,w,t}RV_t^{2W} + \beta_{2W,q,t}RV_{t-6:t}^{2W} + \beta_{2W,y,t}RV_{t-26:t}^{2W} \\
&\quad + \sum_{k=0}^{n-1} \sum_{m=1}^n \eta_{2W,m,t-k}x_{m,t-k} + \epsilon_{2W,t+1} \\
\widehat{RV}_{t+1}^M &= c_M + \beta_{M,W,t}RV_t^M + \beta_{M,q,t}RV_{t-3:t}^M + \beta_{M,y,t}RV_{t-12:t}^M \\
&\quad + \sum_{k=0}^n \eta_{m=1}^{m,t-k}x_{m,t-k} + \epsilon_{M,t+1}
\end{aligned}$$

The time-varying coefficients can be modeled using various methodologies, such as the Markov Chain Monte Carlo (MCMC) method, the Kalman filter method, and the kernel smoothing method. In this research, we apply the kernel smoothing method to estimate the parameters.

### Least Absolute Shrinkage and Selection Operator

Least Absolute Shrinkage and Selection Operator (LASSO) reduces the risk of overfitting by adding a penalty term to the cost function, and it performs variable selection:

$$\begin{aligned}
&\min_{\phi, \eta} \left\{ \left( RV_{t+1} - \left( \sum_{i=0}^{w-1} \phi_i RV_{t-i} + \sum_{k=0}^{w-1} \sum_{m=1}^n \eta_{m,t-k} x_{m,t-k} \right) \right)^2 \right\}, \\
&\text{subject to } \sum |\phi| + \sum |\eta| < c \text{ (constant)}.
\end{aligned}$$

It also has a useful application in the field of energy forecasting, such as research on the oil price predictability in Miao et al. (2017) and Ma et al. (2018).

### Elastic Net

Elastic Net is the combinatorial extension of the Ridge and the Lasso. Although the Ridge and the Lasso enhance the prediction accuracy with the bias-variance tradeoff, they have different performance advantages depending on the data. The Lasso is better when there are fewer significant variables, and the Ridge is better when there are vice versa. The Elastic Net considers both penalty terms of the Ridge and the Lasso. Because Elastic Net has both  $L1$  and  $L2$  norms, it can perform well with large data.

### Decision Tree Regression

Decision Tree Regression (DTR) builds models in the form of a tree structure by dividing the data set into smaller subsets while gradually elaborating related decision trees. The predicted outcome becomes continuous real values.

Hastie et al. (2009) formally describes this concept as follows: when  $X$  is an input variable and defining half-planes as  $R_1(j, s) = \{X|X_j \leq s\}$  and  $R_2(j, s) = \{X|X_j > s\}$ , DTR tries to find the splitting variable  $j$  and split point  $s$  that solve

$$\min_{j,s} \left[ \min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right].$$

where  $y_i$ 's are true values. After finding the best split, repeat this process for each split region. They further discuss some important questions like how large a should tree grow. See Hastie et al. (2009) for the full detailed explanation.

### Random Forest Regression

Random Forest Regression (RFR) is the learning algorithm that uses an ensemble learning method for several decision trees. Even if trees can make an excellent predictive model, they can be very noisy. RFR generates multiple decision trees and averages them to reduce noise and the risk of overfitting. See Hastie et al. (2009) for detailed descriptions.

### Gradient Boosting Regression

Gradient Boosting Regression (GBR) produces a model from an ensemble of weak predictive models. One obvious candidate for the weak predictor will be the tree. The

gradient boosting algorithm updates its predictor by calculating the negative gradient of the loss criteria, then regress a tree to those residuals. A detailed explanation for the algorithm is in Hastie et al. (2009).

### Artificial Neural Network

Artificial Neural Network (ANN) is a computational model inspired by the human brain, designed to recognize patterns and make decisions. It consists of interconnected layers of nodes (neurons), including an input layer, one or more hidden layers, and an output layer. Each connection has a weight that is adjusted during training, allowing the network to learn from data. See Goodfellow (2016) for more details.

### Recurrent Neural Network

Recurrent Neural Network (RNN) is a specific type of artificial neural network designed to process sequential data. Unlike feedforward networks, RNNs have connections that form directed cycles, allowing them to retain information from earlier inputs in the sequence. This makes RNNs particularly effective for tasks like time series prediction. For more details, see Goodfellow (2016).

Using all input data to train models often raises several caveats. The most prominent issue will be training time. As the number of explanatory variables rises, models suffer from a surge in training time. Also, unnecessary variables could hurt the out-of-sample performance. To prevent this, we select features for each sample window before we start training the model.<sup>9</sup> During training, we use the grid search method to find optimal hyperparameters.

For feature selection, we choose the  $f$ -regression method. It first computes the correlation between the dependent variable and each factor. Then it is converted to  $F$ -score, and  $p$ -value of  $F$ -score (Pedregosa et al. (2011)). The linearity of the  $f$ -regression method can quickly calculate the effect of each factor on dependent variables, which is adequate in this paper that uses multiple factors. After calculating  $F$ -scores and  $p$ -values, we

---

<sup>9</sup>Observing the high performance of the HAR-X model, we make an educated guess that the lagged volatility and mean of volatility have high explanatory power. These variables are exempt from the feature selection stage and always belong to the set of explanatory variables.

choose a defined number of factors from the highest score. We always include lagged variables as explanatory variables even if the result of the f-regression method rejects their efficiency to make the machine learning models more coherent with HAR-X.

## 4. Empirical Analysis

### 4.1. *The Measure of Evaluation and Robustness*

Although the  $R^2$  measure is more widely accepted in linear models, it is one of the most intuitive measures for evaluating the performance of a model. In this paper, we examine two  $R^2$  measures:  $\hat{R}_{oos}^2$  score and our version of the out-of-sample  $R_{oos}^2$  score. Let  $y_t$ ,  $\bar{y}_t$ ,  $\hat{y}_t$  and  $\hat{y}_{Bench,t}$  denote true values, the average of true values, predicted values, and the predicted values of the benchmark model HAR, respectively. First, conventionally,  $R_{oos}^2$  score is defined as:

$$R_{oos}^2 = 1 - \frac{\sum_t (y_t - \hat{y}_t)^2}{\sum_t (y_t - \bar{y}_t)^2}$$

Numerous literature argue that  $\bar{y}_t$  should be replaced by the benchmark forecast,  $\hat{y}_{Bench,t}$  (e.g., Christiansen et al. (2012), Wang et al. (2018)). Although the historical average is powerful for forecasting asset returns, it fails to predict asset return volatility. However, the benchmark model HAR has shown a slight decline in its forecasting performance compared to other models recently. So, we introduce a more strict performance measure,  $\hat{R}_{oos}^2$  score:

$$\hat{R}_{oos}^2 = 1 - \frac{\sum_t (y_t - \hat{y}_t)^2}{\sum_t (y_t - \hat{y}_{Bench,t})^2}$$

Second, we define  $\bar{y}_t$  in  $\hat{R}_{oos}^2$  score as the average value of future true values. Note that  $\bar{y}_t$  is not just the historical value, but the average value of our future target variables. Predicted values are stacked outputs from rolling out-of-sample prediction described in Section 3. Therefore,  $R_{oos}^2$  can compare the real-time forecasting performances of various models. Furthermore, the subperiod  $R_{oos}^2$  score is calculated to check the robustness.  $R_{oos}^2$  in the recent 10, 5, 3, 2, and 1-year intervals are measured ( $R_{oos,t}^2$ ,  $t \in \{10, 5, 3, 2, 1\}$ ),

respectively). Furthermore, the stationarity of the residual  $\hat{y}_t - y_t$  is guaranteed in most cases. It makes our model economically meaningful and safe to use performance metrics.<sup>10</sup> As additional metrics, we practice the Model Confidence Set (MCS) introduced by Hansen et al. (2011). For each time horizon - weekly, bi-weekly, and monthly - a set of forecasting models is constructed with a given level of confidence of 0.05. The MCS shows that the best accurate models are different for each time horizon.

In order to test the statistical significance, the Diebold-Mariano (DM) test is employed.<sup>11</sup> The results indicate that the majority of high-performing models exhibit statistically significant differences from the benchmark. We employ the stationarity of the residuals as a proxy for the statistical significance of machine learning models, as these models are unable to provide statistical significance in the form of  $F$ - or  $p$ - values. The majority of the models satisfy the 5% criterion, however, as the periods increase, the value tends to increase (although in most cases, it remains below 5%).

---

<sup>10</sup>We also tested the predictability of residuals  $\hat{y}_t - y_t$  using our models. However, none of the models could predict the one-step-ahead residual with a positive  $R^2_{oss}$ .

<sup>11</sup>More details about DM test is in Appendix C

Forecasting Model	Number of Feature Selection	$\hat{R}_{oos}^2$					
		With Uncertainty			Without Uncertainty		
		M	2W	W	M	2W	W
HAR	-	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)
TV-HAR	-	-0.0203 (2.079E-13)	0.2508 (2.648E-19)	0.2392* (1.421E-19)	-0.0203 (2.079E-13)	0.2508 (2.648E-19)	0.2392* (1.421E-19)
HAR-X	0	-0.0203 (2.079E-13)	0.2508 (2.648E-19)	0.2392* (1.421E-19)	-0.0203 (2.079E-13)	0.2508 (2.648E-19)	0.2392* (1.421E-19)
	10	0.0958 (4.276E-11)	0.3445 (8.056E-16)	0.2528* (1.967E-17)	-0.2564 (2.145E-07)	0.2779 (1.239E-07)	0.1377* (3.584E-14)
	20	0.1260 (4.334E-12)	0.3085 (2.286E-15)	0.2505* (5.876E-17)	0.0739 (1.379E-15)	0.2377 (3.635E-17)	0.0311 (4.732E-15)
	all	<b>0.3708*</b> (4.671E-11)	<b>0.4154*</b> (3.443E-13)	<b>0.3236*</b> (3.356E-12)	0.0125 (6.184E-17)	0.2106 (1.366E-16)	-1.2143 (6.322E-06)
LASSO	0	0.0599 (5.319E-15)	0.3018 (4.384E-20)	0.2756* (1.552E-16)	0.0599 (5.319E-15)	0.3018* (4.384E-20)	0.2756* (1.552E-16)
	10	0.1964 (1.983E-14)	0.3622* (1.719E-19)	0.2766* (1.717E-16)	0.0452 (1.151E-16)	0.2335 (2.426E-15)	0.2605* (2.258E-17)
	20	0.2110 (1.528E-14)	0.3623* (1.723E-19)	0.2766* (1.718E-16)	0.0550 (6.858E-17)	0.2458 (3.340E-15)	0.2603* (2.308E-17)
	all	<b>0.4007</b> (2.232E-17)	<b>0.4512*</b> (8.323E-14)	<b>0.3547*</b> (5.456E-14)	0.2554* (7.102E-17)	0.3980* (5.951E-16)	0.2439* (1.809E-14)
ElasticNet	0	0.0613 (6.682E-15)	0.2985 (3.345E-20)	0.2698* (1.053E-17)	0.0613 (6.682E-15)	0.2985* (3.345E-20)	0.2698* (1.053E-17)
	10	0.2137 (2.232E-14)	0.3849* (4.921E-19)	0.2969* (2.775E-17)	0.1537* (3.724E-15)	0.3245* (3.286E-14)	0.2714* (3.825E-18)
	20	0.2237 (2.426E-14)	0.3807* (3.516E-19)	0.2971* (3.317E-17)	0.2023* (1.340E-20)	0.3319* (1.301E-20)	0.2726* (5.200E-18)
	all	<b>0.4043*</b> (3.434E-17)	<b>0.4397*</b> (3.443E-14)	<b>0.3533*</b> (5.409E-14)	0.2396* (6.246E-17)	0.3884* (8.830E-17)	0.2007* (1.019E-13)
DTR	0	0.0513 (1.209E-12)	-0.0220 (7.737E-08)	-0.2066 (1.258E-17)	0.0240 (2.446E-12)	-0.1064 (2.973E-08)	-0.1085 (7.020E-16)
	10	0.3294* (4.235E-14)	<b>0.3065</b> (3.638E-12)	-0.2111 (7.734E-15)	0.2264* (7.963E-16)	0.0022 (4.163E-15)	-0.2102 (1.461E-19)
	20	<b>0.3452*</b> (2.868E-15)	0.2293 (5.137E-13)	-0.2178 (4.559E-18)	0.2755* (3.774E-14)	-0.0087 (1.436E-13)	<b>0.0697</b> (1.118E-12)
	all	0.3308* (1.535E-15)	0.1281 (1.858E-11)	-0.2805 (7.418E-19)	0.0503 (5.726E-08)	0.1752 (2.422E-18)	-0.0473 (3.920E-13)
RFR	0	0.2827* (1.293E-12)	0.3737* (2.818E-16)	0.2297 (1.783E-17)	0.2823* (1.240E-12)	0.3720* (2.061E-16)	0.2428* (3.330E-17)
	10	0.4623* (3.638E-16)	<b>0.4884*</b> (1.350E-15)	0.3434* (3.149E-16)	0.4228* (1.145E-13)	0.4060* (2.940E-18)	0.2760* (2.131E-15)
	20	<b>0.4909*</b> (1.028E-15)	0.4678* (2.964E-17)	0.3743* (2.684E-17)	0.4216* (1.606E-13)	0.4604* (3.747E-17)	0.3352* (2.574E-16)
	all	0.4843* (3.575E-15)	0.4668* (2.657E-17)	<b>0.3751*</b> (1.246E-17)	0.4175* (8.148E-14)	0.4470* (3.617E-18)	0.3328* (9.301E-17)
GBR	0	0.2160 (5.184E-13)	0.3269 (3.493E-17)	0.2092 (6.703E-16)	0.2296* (6.920E-13)	0.3205* (2.413E-17)	0.2233* (3.080E-16)
	10	0.4292* (1.710E-15)	<b>0.4891*</b> (1.786E-13)	0.3221* (1.266E-15)	0.3102* (1.699E-06)	0.3305* (1.026E-16)	0.0601 (2.916E-19)
	20	<b>0.4783*</b> (2.345E-15)	0.4644* (1.528E-16)	<b>0.3601*</b> (3.625E-16)	0.3382* (2.420E-16)	0.3751* (4.168E-18)	0.2670* (2.752E-19)
	all	0.4641* (1.117E-15)	0.4709* (1.059E-16)	0.3591* (1.366E-16)	0.3575* (3.946E-17)	0.3919* (8.140E-19)	0.2863* (2.269E-17)
ANN	0	-0.0483 (1.238E-13)	0.2008 (8.166E-20)	0.2090 (6.319E-18)	-0.0554 (4.572E-14)	0.2214 (1.274E-19)	0.1738* (4.084E-18)
	10	0.2757* (9.829E-11)	0.4146* (5.968E-13)	<b>0.3610*</b> (3.832E-13)	0.1210 (4.540E-14)	0.4633* (1.464E-13)	0.2837* (8.791E-13)
	20	0.2485 (2.381E-12)	0.3708* (2.239E-15)	0.2801* (4.347E-14)	0.2626* (1.589E-13)	<b>0.5103*</b> (1.832E-12)	0.0482 (4.429E-13)
	all	<b>0.2923*</b> (5.360E-12)	0.3285 (4.946E-16)	0.2618* (1.763E-13)	-2.2483 (2.300E-18)	-0.0852 (3.332E-14)	-0.7148 (1.649E-11)
RNN	0	0.4398* (5.641E-10)	0.3868* (2.328E-15)	0.3121* (8.896E-15)	0.4372* (6.683E-10)	0.3878* (8.813E-16)	0.3132* (1.264E-14)
	10	0.4711* (1.647E-11)	<b>0.4777*</b> (1.101E-13)	0.3956* (7.267E-14)	0.4800* (5.836E-11)	0.4696* (1.699E-12)	<b>0.4070*</b> (6.483E-15)
	20	<b>0.4906*</b> (4.541E-10)	0.4549* (3.328E-15)	0.3841* (1.998E-15)	0.4803* (3.066E-10)	0.4457* (2.522E-07)	0.3488* (3.794E-14)
	all	0.4813* (6.191E-10)	0.4316* (2.055E-13)	0.3347* (4.252E-15)	0.4742* (8.149E-10)	0.4135* (5.288E-14)	0.3280* (4.572E-15)

**Table 2. Forecasting Performance Compared to the Benchmark HAR.**  $\hat{R}_{oos}^2$  is the performance measurement compared to the benchmark HAR. For each model and each forecasting horizon, bold numbers are the best performance measurements with feature selection and data set. For each time horizon, models that pass the MCS test with a significant level of 0.05 have an asterisk with each performance measure. The numbers in parentheses indicate the  $p$ -values of residual stationarity.

		DM Test					
Forecasting Model	Number of Feature Selection	With Uncertainty			Without Uncertainty		
		M	2W	W	M	2W	W
HAR	-	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)	0.0000 (0.000E+00)
TV-HAR	-	-0.0453 (9.639E-01)	1.0316 (3.023E-01)	1.3210 (1.865E-01)	-0.0453 (9.639E-01)	1.0316 (3.023E-01)	1.3210 (1.865E-01)
HAR-X	0	-0.0453 (9.639E-01)	1.0316 (3.023E-01)	1.3210 (1.865E-01)	-0.0453 (9.639E-01)	1.0316 (3.023E-01)	1.3210 (1.865E-01)
	10	0.2264 (8.209E-01)	1.1966 (2.315E-01)	1.2413 (2.145E-01)	-0.4538 (6.500E-01)	0.8356 (4.034E-01)	0.6654 (5.058E-01)
	20	0.2949 (7.680E-01)	1.1187 (2.633E-01)	1.1772 (2.391E-01)	0.1583 (8.742E-01)	0.7380 (4.605E-01)	0.1904 (8.490E-01)
	all	<b>0.7735</b> (4.392E-01)	<b>1.3406</b> (1.801E-01)	<b>1.6385</b> (1.013E-01)	0.0251 (9.800E-01)	0.6020 (5.471E-01)	-2.1993 (2.786E-02)
LASSO	0	0.1370 (8.910E-01)	1.2562 (2.091E-01)	1.5992 (1.098E-01)	0.1370 (8.910E-01)	1.2562 (2.091E-01)	1.5992 (1.098E-01)
	10	0.4524 (6.510E-01)	1.4092 (1.588E-01)	1.6055 (1.084E-01)	0.0950 (9.243E-01)	0.9466 (3.439E-01)	1.5345 (1.249E-01)
	20	0.4886 (6.251E-01)	1.4095 (1.587E-01)	1.6054 (1.084E-01)	0.1163 (9.074E-01)	1.0021 (3.163E-01)	1.5334 (1.252E-01)
	all	<b>0.8641</b> (3.875E-01)	<b>1.5312</b> (1.257E-01)	<b>1.9618</b> (4.979E-02)	0.5400 (5.892E-01)	1.3853 (1.659E-01)	1.3523 (1.763E-01)
ElasticNet	0	0.1396 (8.890E-01)	1.2414 (2.145E-01)	1.5810 (1.139E-01)	0.1396 (8.890E-01)	1.2414 (2.145E-01)	1.5810 (1.139E-01)
	10	0.4955 (6.203E-01)	1.4632 (1.434E-01)	1.7104 (8.719E-02)	0.3447 (7.303E-01)	1.2239 (2.210E-01)	1.5775 (1.147E-01)
	20	0.5184 (6.042E-01)	1.4592 (1.445E-01)	1.7080 (8.764E-02)	0.4584 (6.466E-01)	1.2538 (2.099E-01)	1.5816 (1.137E-01)
	all	<b>0.8611</b> (3.892E-01)	<b>1.5033</b> (1.328E-01)	<b>1.9617</b> (4.980E-02)	0.5064 (6.126E-01)	1.3055 (1.917E-01)	1.0979 (2.723E-01)
DTR	0	0.1157 (9.079E-01)	-0.0534 (9.574E-01)	-0.7162 (4.739E-01)	0.0542 (9.568E-01)	-0.2489 (8.034E-01)	-0.4847 (6.279E-01)
	10	0.6236 (5.329E-01)	<b>0.9530</b> (3.406E-01)	-0.6466 (5.179E-01)	0.4396 (6.602E-01)	0.0071 (9.944E-01)	-0.5948 (5.520E-01)
	20	<b>0.6709</b> (5.023E-01)	0.7310 (4.648E-01)	-0.6852 (4.932E-01)	0.5042 (6.141E-01)	-0.0303 (9.759E-01)	<b>0.3664</b> (7.141E-01)
	all	0.6438 (5.197E-01)	0.3978 (6.908E-01)	-0.8483 (3.963E-01)	0.0800 (9.362E-01)	0.5657 (5.716E-01)	-0.2662 (7.901E-01)
RFR	0	0.5910 (5.545E-01)	1.2729 (2.031E-01)	1.2499 (2.113E-01)	0.5902 (5.551E-01)	1.2656 (2.057E-01)	1.3154 (1.884E-01)
	10	0.9039 (3.660E-01)	<b>1.6849</b> (9.201E-02)	1.9719 (4.862E-02)	0.8215 (4.113E-01)	1.5325 (1.254E-01)	1.4702 (1.415E-01)
	20	<b>0.9639</b> (3.351E-01)	1.6549 (9.795E-02)	2.0061 (4.484E-02)	0.8104 (4.177E-01)	1.6460 (9.976E-02)	1.7020 (8.875E-02)
	all	0.9434 (3.455E-01)	1.6520 (9.854E-02)	<b>2.0083</b> (4.461E-02)	0.8075 (4.194E-01)	1.6108 (1.072E-01)	1.7388 (8.207E-02)
GBR	0	0.4689 (6.391E-01)	1.0528 (2.924E-01)	1.1469 (2.514E-01)	0.4946 (6.209E-01)	1.0299 (3.031E-01)	1.2093 (2.265E-01)
	10	0.8215 (4.114E-01)	<b>1.6461</b> (9.975E-02)	1.8191 (6.889E-02)	0.5836 (5.595E-01)	1.2225 (2.215E-01)	0.2408 (8.097E-01)
	20	<b>0.9446</b> (3.449E-01)	1.5738 (1.155E-01)	<b>1.8751</b> (6.078E-02)	0.6352 (5.253E-01)	1.4954 (1.348E-01)	1.4490 (1.473E-01)
	all	0.9041 (3.659E-01)	1.6231 (1.046E-01)	<b>1.9438</b> (5.192E-02)	0.6768 (4.986E-01)	1.5035 (1.327E-01)	1.5062 (1.320E-01)
ANN	0	-0.1063 (9.153E-01)	0.8240 (4.100E-01)	1.2178 (2.233E-01)	-0.1211 (9.036E-01)	0.9147 (3.604E-01)	1.0322 (3.020E-01)
	10	0.6114 (5.409E-01)	1.2917 (1.965E-01)	<b>2.1024</b> (3.552E-02)	0.2777 (7.813E-01)	1.2259 (2.202E-01)	1.8050 (7.107E-02)
	20	0.5951 (5.518E-01)	1.2480 (2.120E-01)	1.3910 (1.642E-01)	0.5359 (5.920E-01)	<b>1.3297</b> (1.836E-01)	0.2697 (7.874E-01)
	all	<b>0.6462</b> (5.182E-01)	1.1663 (2.435E-01)	1.3365 (1.814E-01)	-0.9776 (3.283E-01)	-0.2081 (8.351E-01)	-1.1137 (2.654E-01)
RNN	0	0.8373 (4.024E-01)	1.3570 (1.748E-01)	1.8215 (6.852E-02)	0.8281 (4.076E-01)	1.3550 (1.754E-01)	1.8357 (6.641E-02)
	10	0.9064 (3.648E-01)	<b>1.6585</b> (9.721E-02)	<b>2.2714</b> (2.313E-02)	0.9412 (3.466E-01)	1.6314 (1.028E-01)	<b>2.2082</b> (2.723E-02)
	20	<b>0.9504</b> (3.419E-01)	1.5747 (1.153E-01)	2.2383 (2.520E-02)	0.9256 (3.547E-01)	1.5598 (1.188E-01)	2.1720 (2.985E-02)
	all	0.9364 (3.490E-01)	1.5226 (1.279E-01)	2.0474 (4.062E-02)	0.9155 (3.599E-01)	1.4484 (1.475E-01)	2.0146 (4.394E-02)

**Table 3. Additional Metrics.** The DM contain  $p$ -values for whether each model has significantly different accuracy compared to benchmark HAR. The numbers in parentheses indicate the  $p$ -values of residual stationarity.



## 4.2. Empirical Results and Potential Explanation

**Overall Performance** The out-of-sample test provides evidence that machine learning models are capable of achieving high levels of accuracy in forecasting oil volatility. This is evidenced by a notable enhancement in the adjusted R-squared ( $R^2$ ) scores relative to the benchmark HAR model. The most commonly used time-series models, such as HAR-X and TV-HAR, demonstrate a less impressive performance than machine learning models, yet they also exhibit a superior performance compared to the HAR. The results inhibit that linear or simple machine learning models, including LASSO, ElasticNet, and DTR, demonstrate inconsistent performance. These models typically demonstrate superior performance when forecasting past periods, whereas their efficacy tends to diminish in current periods. Also, the number of feature selection plays an important role in performance. Nonlinear machine learning models with many small predictors such as the RFR and GBR can reach an impressive out-of-sample  $R^2$  score.

Furthermore, the number of feature selections is a substantial determinant of performance. The performance of linear models, including HAR-X, LASSO, and ElasticNet, is enhanced as the number of selected features increases. Conversely, nonlinear machine learning models, such as RFR and GBR, demonstrated an exceptional out-of-sample  $R^2$  score with 20 features. This suggests that an excess of features may negatively impact the performance of the models. Deep learning models, including ANN and RNN, yield varying results. Given the sensitivity of deep learning models to hyperparameters and their tendency to require more training, the outcomes of deep learning models may be inconsistent and subject to change over time. However, the performance of RNN is one of the most notable models in recent times.

As described in section 4.1, we provide a wide range of additional analysis and robustness checks. The MCS test and the Diebold-Mariano test verifies which models comparatively outperform the benchmark and other models. We also provide the robustness checks of the forecast horizon, the number of feature selections, the sub-period length of performance measurement, and an additional set of uncertainty indices.

The MCS test in Table 2 shows that the best-performing models vary over time horizons and a set of input data. For the monthly forecast, the random forest and RNN models work well in the out-of-sample. However, for the weekly forecast, some

of the models even underperform the TV-HAR, which implies the strong short-term momentum of the weekly realized volatility.

Table 3 provides additional metrics - the DM test. We present the statistic and  $p$ -values of the DM test with the benchmark model. Although the DM test is another comparative statistic about forecasting performance, the results can be different from those of the MCS test. This is because the MCS test is a statistical test for multiple models, excluding inferior models well, while the DM test is the one-to-one comparison with the benchmark. Most models from HAR-X to RNN show positive statistics, which means that the models are better than the benchmark HAR. Interestingly, some models such as DTR, GBR, and ANN are not differentiated statistically with HAR. One explanation could be lagged variables. The two most important factors in the HAR model are short-term and long-term lagged variables. In DTR and GBR, the lagged variables could affect DTR and GBR the most, and therefore, not statistically differentiated from the HAR model.

Forecasting Model	Number of Feature Selection	$R^2_{OOS}$	$R^2_{OOS,10}$	$R^2_{OOS,5}$	$R^2_{OOS,3}$	$R^2_{OOS,2}$	$R^2_{OOS,1}$
HAR	-	-0.8224 (2.273E-03)	-0.8918 (4.362E-03)	-1.0895 (2.453E-01)	-0.1665 (9.286E-01)	-0.1312 (1.728E-02)	-1.0531 (7.422E-02)
TV-HAR	-	-0.1154 (5.893E-21)	-0.1759 (4.605E-19)	-0.2851 (1.943E-05)	-0.1238 (4.493E-01)	0.3608 (2.088E-01)	0.1883 (3.401E-02)
HAR-X	0	-0.8224 (2.273E-03)	-0.8918 (4.362E-03)	-1.0895 (2.453E-01)	-0.1665 (9.286E-01)	-0.1312 (1.728E-02)	-1.0531 (7.422E-02)
	10	-0.1838 (5.802E-05)	-0.2409 (2.546E-04)	-0.3486 (3.216E-02)	-0.9011 (1.565E-02)	-1.0257 (3.521E-04)	-0.2328 (1.051E-02)
	20	0.2034 (3.729E-07)	0.1826 (1.876E-06)	0.1161 (1.012E-01)	-1.4794 (6.567E-04)	-1.6086 (3.814E-03)	-0.8476 (1.001E-03)
	all	0.0591 (9.227E-08)	0.0706 (7.447E-07)	0.0697 (7.134E-09)	-2.0926 (1.938E-01)	-1.7339 (2.477E-01)	-0.4224 (9.975E-01)
		-0.0256 (5.237E-11)	-0.0788 (6.249E-10)	-0.1795 (4.716E-05)	-0.0518 (5.148E-01)	0.3053 (8.338E-01)	-0.0484 (7.259E-02)
LASSO	0	-0.1838 (5.802E-05)	-0.2409 (2.546E-04)	-0.3486 (3.216E-02)	-0.9011 (1.565E-02)	-1.0257 (3.521E-04)	-0.2328 (1.051E-02)
	10	0.1452 (2.770E-16)	0.1024 (1.519E-14)	0.0301 (2.777E-02)	-0.0758 (1.269E-10)	0.3622 (1.406E-04)	0.2709 (3.048E-02)
	20	0.1591 (5.892E-16)	0.1170 (3.173E-14)	0.0353 (2.446E-05)	-0.0854 (5.776E-09)	0.3333 (1.099E-04)	0.2862 (3.157E-02)
	all	<b>0.4333</b> (2.474E-17)	<b>0.4293</b> (2.084E-16)	<b>0.3940</b> (2.004E-09)	-0.4985 (2.058E-02)	-0.1757 (2.772E-01)	0.2427 (5.396E-01)
		-0.0376 (1.311E-10)	-0.0911 (1.308E-09)	-0.1924 (8.332E-05)	-0.0595 (5.239E-01)	0.3015 (7.996E-01)	-0.0452 (6.511E-02)
ElasticNet	0	0.2946 (1.414E-15)	0.2612 (5.491E-14)	0.2136 (1.169E-05)	0.0662 (1.267E-01)	0.4722 (5.872E-05)	0.3081 (1.638E-02)
	10	0.3418 (1.298E-13)	0.3113 (2.513E-12)	0.2602 (1.372E-06)	-0.0608 (3.984E-02)	0.4135 (4.591E-01)	0.4078 (5.725E-01)
	20	0.4092 (1.381E-16)	0.4075 (9.441E-16)	0.3714 (4.328E-03)	-0.5334 (1.420E-02)	-0.2036 (2.680E-01)	0.1589 (6.463E-01)
	all	-0.2242 (5.670E-19)	-0.2852 (2.993E-10)	-0.3480 (8.087E-05)	-2.5173 (3.350E-06)	-4.5398 (9.969E-01)	-9.0861 (3.418E-02)
		0.2721 (4.367E-15)	0.2786 (9.592E-14)	0.3157 (6.653E-04)	0.1458 (6.463E-07)	-0.1037 (1.000E+00)	-0.4031 (9.925E-01)
DTR	0	-0.0420 (1.149E-15)	-0.0536 (5.663E-14)	-0.0797 (1.362E-09)	0.0129 (2.201E-06)	0.0686 (1.621E-03)	0.1317 (1.744E-01)
	10	-0.3388 (1.001E-14)	-0.3719 (5.967E-18)	-0.4815 (1.616E-05)	-0.1227 (3.519E-06)	-0.0804 (1.485E-01)	-0.4701 (9.981E-01)
	20	0.1644 (6.369E-17)	0.1218 (2.867E-15)	0.0765 (6.978E-04)	-0.5679 (6.444E-09)	-0.2707 (6.657E-03)	-1.2435 (9.910E-01)
	all	0.3904 (2.400E-16)	0.3652 (8.450E-15)	0.3311 (4.116E-02)	<b>0.3023</b> (3.036E-07)	0.5334 (4.232E-04)	0.4312 (1.389E-01)
		0.3509 (6.550E-17)	0.3233 (2.254E-15)	0.2802 (7.894E-02)	0.1156 (5.473E-07)	0.4165 (4.323E-01)	0.4939 (9.963E-01)
RFR	0	0.3366 (4.679E-17)	0.3092 (1.436E-15)	0.2620 (8.484E-02)	0.1012 (1.527E-01)	0.4069 (8.992E-01)	0.4465 (9.864E-01)
	10	0.1262 (2.970E-17)	0.0821 (1.418E-15)	0.0368 (3.394E-04)	-0.5351 (6.945E-08)	-0.3546 (3.171E-02)	-1.1792 (1.000E+00)
	20	0.2712 (1.077E-14)	0.2489 (2.654E-13)	0.2077 (1.503E-02)	0.1331 (3.247E-01)	0.2344 (8.151E-05)	0.2781 (2.180E-01)
	all	0.1082 (1.293E-19)	0.0739 (6.170E-18)	0.0023 (2.454E-02)	-0.4301 (2.570E-06)	-0.6828 (8.497E-01)	0.3736 (9.956E-01)
		0.1630 (1.314E-19)	0.1266 (4.881E-18)	0.0598 (1.564E-02)	-0.1553 (2.266E-06)	-0.0990 (1.317E-08)	0.3394 (5.420E-02)
GBR	0	-0.1792 (5.113E-12)	-0.2224 (4.647E-12)	-0.3356 (5.583E-05)	-0.0549 (7.479E-01)	0.2994 (5.361E-01)	0.1402 (6.861E-02)
	10	0.3931 (4.965E-06)	0.3763 (2.072E-05)	0.3590 (1.444E-06)	-0.1711 (2.086E-01)	0.1225 (9.991E-01)	-0.0981 (4.872E-01)
	20	0.4052 (2.865E-16)	0.3891 (2.758E-15)	0.3680 (6.164E-06)	-0.6005 (7.019E-06)	-0.6262 (2.046E-03)	-0.8552 (1.657E-01)
	all	-3.7435 (4.146E-18)	-3.9785 (1.864E-16)	-4.6024 (3.584E-09)	-0.9025 (5.994E-07)	-0.3906 (7.837E-02)	-0.8096 (1.790E-01)
		0.0488 (5.805E-10)	0.0102 (4.372E-09)	-0.0347 (8.138E-04)	-0.0255 (7.119E-01)	0.1527 (2.079E-01)	-0.4145 (3.593E-01)
ANN	0	0.0892 (1.141E-10)	0.0449 (1.593E-09)	-0.0159 (3.448E-04)	0.1064 (6.350E-08)	0.5063 (6.999E-03)	0.2701 (1.280E-01)
	10	0.1032 (6.659E-10)	0.0676 (7.596E-09)	0.0215 (1.424E-03)	0.1578 (2.054E-01)	<b>0.5888</b> (5.465E-05)	<b>0.6056</b> (3.197E-02)
	20	0.0663 (6.894E-10)	0.0294 (2.572E-09)	0.0061 (2.506E-04)	0.2718 (3.152E-07)	0.5737 (2.702E-02)	0.5002 (9.885E-01)
	all						

**Table 4. Monthly Forecasting Performance.**  $R^2_{OOS}$  scores for entire period and subperiod ( $t \in \{10, 5, 3, 2, 1\}$ ) are represented. The numbers in parentheses indicate the  $p$ -values of residual stationarity. The bold numbers are the best performance measurements for each subperiod.

Forecasting Model	Number of Feature Selection	$R^2_{OOS}$	$R^2_{OOS,10}$	$R^2_{OOS,5}$	$R^2_{OOS,3}$	$R^2_{OOS,2}$	$R^2_{OOS,1}$
HAR	-	-0.4331 (1.471E-15)	-0.5165 (2.663E-12)	-0.6729 (2.883E-07)	0.0788 (3.739E-12)	0.0702 (2.547E-03)	-0.0623 (2.408E-06)
TV-HAR	-	0.1240 (8.207E-18)	0.0891 (7.085E-14)	0.0197 (1.216E-07)	0.2160 (1.090E-12)	0.2571 (1.267E-12)	0.1236 (6.754E-06)
HAR-X	0	-0.4331 (1.471E-15)	-0.5165 (2.663E-12)	-0.6729 (2.883E-07)	0.0788 (3.739E-12)	0.0702 (2.547E-03)	-0.0623 (2.408E-06)
	10	-0.0442 (1.502E-05)	-0.0670 (3.839E-10)	-0.1590 (2.935E-26)	-0.1486 (2.547E-13)	-0.1498 (2.798E-13)	0.0735 (1.778E-09)
	20	-0.0568 (9.822E-14)	-0.0768 (8.450E-11)	-0.1621 (1.794E-05)	-0.6356 (6.374E-14)	-0.9367 (2.520E-12)	-1.9959 (2.046E-04)
	all	-0.4293 (1.761E-13)	-0.1344 (1.518E-09)	-0.2318 (1.534E-05)	-1.1232 (6.778E-13)	-1.8694 (6.827E-10)	-3.6012 (2.259E-03)
LASSO	0	0.1724 (3.072E-17)	0.1455 (4.827E-14)	0.0814 (1.725E-08)	0.2501 (1.690E-12)	0.2458 (8.294E-12)	0.0157 (1.583E-05)
	10	0.0692 (3.391E-17)	0.0312 (6.652E-14)	-0.0522 (1.655E-08)	0.2273 (1.077E-10)	0.2668 (1.866E-10)	0.0754 (5.149E-05)
	20	0.0871 (1.555E-17)	0.0511 (3.371E-14)	-0.0289 (1.464E-08)	0.1833 (1.917E-11)	0.2666 (2.988E-12)	0.0814 (3.661E-05)
	all	0.0867 (1.900E-12)	0.1744 (5.250E-10)	0.0995 (1.554E-05)	-0.5802 (2.665E-12)	-1.2603 (1.210E-11)	-1.6866 (9.654E-03)
ElasticNet	0	0.1713 (2.450E-17)	0.1451 (3.507E-14)	0.0810 (1.360E-08)	0.2486 (1.721E-12)	0.2446 (8.137E-12)	0.0197 (1.713E-05)
	10	0.1450 (5.186E-16)	0.1135 (3.881E-13)	0.0423 (1.020E-07)	0.2495 (1.478E-11)	0.2417 (3.881E-12)	0.2120 (9.338E-07)
	20	0.1700 (1.243E-16)	0.1391 (1.359E-13)	0.0715 (7.929E-08)	0.1906 (2.680E-12)	0.2595 (2.098E-14)	0.2136 (4.261E-07)
	all	0.0228 (6.375E-13)	0.1533 (7.425E-11)	0.0768 (2.863E-06)	-0.6958 (2.870E-12)	-1.4011 (4.179E-12)	-2.1249 (4.316E-02)
DTR	0	-0.4325 (8.693E-16)	-0.4692 (2.541E-12)	-0.5355 (4.940E-07)	-3.7725 (3.658E-15)	-1.4038 (7.028E-12)	-2.4123 (1.147E-05)
	10	-0.2340 (2.188E-17)	-0.2539 (1.437E-11)	-0.3301 (8.426E-09)	-1.8767 (1.045E-14)	-0.7251 (1.157E-12)	-1.3815 (1.490E-07)
	20	-0.0839 (0.000E+00)	-0.0952 (2.079E-30)	-0.1367 (5.238E-25)	-1.9534 (2.002E-12)	-0.2968 (2.542E-11)	-0.7026 (1.418E-03)
	all	0.0647 (9.164E-16)	0.0470 (7.741E-13)	0.0123 (1.856E-08)	-0.3377 (2.446E-19)	-0.5570 (6.756E-17)	-0.0361 (2.082E-07)
RFR	0	0.2153 (5.245E-15)	0.2149 (3.728E-12)	0.1811 (9.815E-07)	-0.3692 (2.024E-15)	0.0289 (1.631E-12)	0.0245 (1.216E-01)
	10	0.3123 (2.016E-16)	0.3205 (9.289E-14)	0.2754 (1.423E-08)	0.3384 (1.029E-11)	0.0147 (3.652E-03)	-0.0736 (3.997E-04)
	20	0.4017 (4.199E-15)	0.4056 (4.871E-12)	0.3755 (1.403E-07)	0.4038 (1.540E-12)	0.2424 (4.180E-04)	0.0699 (1.067E-04)
	all	<b>0.4059</b> (1.036E-16)	<b>0.4121</b> (3.783E-14)	<b>0.3833</b> (1.575E-08)	0.3914 (1.180E-12)	0.2325 (2.667E-12)	0.1978 (1.191E-05)
GBR	0	0.1327 (3.405E-15)	0.1364 (3.168E-12)	0.1045 (7.101E-07)	-0.6758 (3.688E-16)	0.0271 (1.539E-12)	-0.1844 (9.086E-02)
	10	0.2331 (5.805E-17)	0.2320 (3.191E-15)	0.1814 (3.698E-09)	0.4332 (3.902E-10)	0.2055 (4.199E-10)	0.1296 (2.195E-05)
	20	0.3058 (1.199E-15)	0.3084 (2.082E-13)	0.2749 (2.461E-08)	0.4065 (9.598E-11)	0.2670 (1.084E-11)	0.1989 (4.096E-05)
	all	0.3687 (4.186E-18)	0.3785 (8.608E-16)	0.3586 (1.349E-10)	0.4026 (4.284E-10)	<b>0.3182</b> (2.447E-04)	0.2105 (2.794E-01)
ANN	0	0.0898 (2.762E-17)	0.0667 (5.953E-14)	-0.0034 (5.195E-09)	0.1598 (6.739E-14)	0.0732 (7.992E-13)	-0.0893 (5.639E-05)
	10	0.0432 (1.201E-09)	0.0096 (1.558E-08)	-0.0860 (2.050E-04)	-0.3812 (4.485E-13)	-1.1930 (7.545E-13)	-3.1175 (8.074E-08)
	20	0.1891 (9.951E-11)	0.1766 (2.494E-30)	0.1195 (8.753E-25)	-0.7372 (1.257E-04)	-0.9783 (8.250E-09)	-1.9248 (1.002E-06)
	all	-1.8634 (4.023E-30)	-2.0457 (9.166E-28)	-2.4374 (4.068E-19)	-0.8875 (1.316E-11)	-0.9708 (3.112E-08)	-1.9250 (1.432E-10)
RNN	0	0.2298 (3.430E-15)	0.2127 (2.047E-12)	0.1725 (8.075E-07)	0.1970 (2.705E-15)	0.1550 (2.888E-13)	0.0850 (7.678E-07)
	10	0.3321 (1.961E-11)	0.3359 (2.146E-09)	0.3025 (3.208E-05)	0.4098 (8.352E-11)	0.2218 (3.023E-11)	0.1539 (3.771E-05)
	20	0.2806 (1.830E-11)	0.2667 (3.291E-09)	0.2188 (2.911E-05)	0.2068 (1.743E-10)	0.1776 (8.148E-13)	0.0173 (2.713E-07)
	all	0.2578 (6.886E-08)	0.2373 (6.112E-06)	0.1996 (2.988E-06)	<b>0.4725</b> (1.072E-10)	0.2964 (1.575E-02)	<b>0.2947</b> (9.153E-05)

**Table 5. Bi-weekly Forecasting Performance.**  $R^2_{oos}$  scores for entire period and subperiod ( $t \in \{10, 5, 3, 2, 1\}$ ) are represented. The numbers in parentheses indicate the  $p$ -values of residual stationarity. The bold numbers are the best performance measurements for each subperiod.

Forecasting Model	Number of Feature Selection	$R^2_{OOS}$	$R^2_{OOS,10}$	$R^2_{OOS,5}$	$R^2_{OOS,3}$	$R^2_{OOS,2}$	$R^2_{OOS,1}$
HAR	-	-0.0714 (6.045E-17)	-0.1217 (9.597E-13)	-0.2127 (2.885E-07)	-0.0059 (6.728E-18)	0.0323 (7.068E-17)	-0.1408 (5.059E-09)
TV-HAR	-	0.3117 (2.151E-19)	0.2990 (6.913E-15)	0.2749 (8.617E-09)	0.1455 (6.334E-27)	0.1225 (8.163E-11)	0.0129 (1.361E-14)
HAR-X	0	-0.0714 (6.045E-17)	-0.1217 (9.597E-13)	-0.2127 (2.885E-07)	-0.0059 (6.728E-18)	0.0323 (7.068E-17)	-0.1408 (5.059E-09)
	10	0.2495 (1.013E-12)	0.2571 (3.066E-10)	0.2279 (2.523E-05)	-0.2038 (2.234E-24)	-0.3675 (1.808E-21)	-0.5944 (7.910E-16)
	20	0.2082 (1.646E-14)	0.2708 (1.611E-11)	0.2487 (6.935E-06)	-0.4295 (6.239E-24)	-0.7749 (1.835E-20)	-1.7198 (3.046E-12)
	all	-1.1823 (4.979E-04)	0.2932 (1.317E-09)	0.2632 (2.614E-19)	-0.4458 (2.301E-22)	-0.9169 (3.547E-17)	-1.8841 (1.661E-08)
LASSO	0	0.3367 (1.174E-14)	0.3321 (1.205E-11)	0.3114 (8.683E-07)	0.1872 (1.618E-26)	0.1673 (2.353E-22)	0.0231 (1.272E-14)
	10	0.3296 (6.190E-16)	0.3239 (1.045E-12)	0.3016 (1.621E-07)	0.1889 (2.387E-26)	0.1700 (3.201E-22)	0.0260 (1.472E-14)
	20	0.3299 (6.322E-16)	0.3239 (1.045E-12)	0.3016 (1.621E-07)	0.1889 (2.387E-26)	0.1700 (3.201E-22)	0.0260 (1.472E-14)
	all	0.2928 (2.227E-13)	0.3658 (1.310E-10)	0.3338 (3.810E-06)	-0.1595 (3.262E-21)	-0.4780 (7.187E-17)	-0.7846 (3.602E-10)
ElasticNet	0	0.3398 (1.618E-16)	0.3332 (5.191E-13)	0.3124 (1.502E-07)	0.1852 (1.320E-26)	0.1645 (1.758E-22)	0.0248 (1.237E-14)
	10	0.3577 (1.943E-17)	0.3509 (1.161E-13)	0.3335 (4.603E-08)	0.2113 (3.535E-26)	0.1847 (2.769E-22)	0.0572 (1.350E-14)
	20	0.3584 (2.920E-17)	0.3521 (1.258E-13)	0.3350 (4.836E-08)	0.1975 (2.361E-26)	0.1786 (5.986E-23)	<b>0.0735</b> (4.284E-15)
	all	0.2568 (2.294E-12)	0.3645 (4.002E-10)	0.3332 (7.815E-06)	-0.2146 (5.368E-22)	-0.5612 (1.015E-17)	-1.0275 (4.570E-10)
DTR	0	-0.2206 (2.153E-16)	-0.1725 (1.191E-11)	-0.1309 (7.171E-06)	-1.8657 (3.911E-22)	-1.4720 (9.982E-19)	-2.3459 (5.843E-09)
	10	-0.2822 (1.195E-19)	-0.2305 (1.764E-15)	-0.2959 (1.667E-08)	-5.9182 (1.224E-25)	-13.7236 (1.383E-20)	-41.8268 (2.284E-12)
	20	-0.0511 (7.846E-20)	0.0016 (2.562E-08)	-0.0174 (1.872E-18)	-0.9437 (2.031E-25)	-2.4880 (1.856E-20)	-5.3090 (1.641E-11)
	all	-0.1687 (5.427E-14)	-0.1334 (1.110E-10)	-0.1704 (2.366E-07)	-1.9978 (2.579E-26)	-4.7846 (1.276E-20)	-12.5673 (1.110E-12)
RFR	0	0.2462 (1.157E-15)	0.2516 (3.219E-12)	0.2453 (3.336E-06)	0.0239 (6.389E-25)	0.0903 (6.060E-10)	-0.0627 (4.391E-09)
	10	0.3319 (2.497E-14)	0.3463 (1.239E-15)	0.3125 (2.334E-08)	-0.6355 (5.984E-25)	-2.0051 (9.341E-20)	-6.5592 (2.090E-12)
	20	0.3900 (1.917E-15)	0.4037 (4.770E-12)	0.3846 (1.325E-06)	-0.3368 (4.639E-27)	-1.3978 (1.119E-14)	-5.0238 (9.976E-09)
	all	0.3877 (4.186E-16)	0.3990 (9.884E-13)	0.3768 (2.392E-07)	-0.2965 (1.884E-26)	-1.2702 (1.573E-14)	-4.4404 (1.375E-13)
GBR	0	0.2203 (2.421E-14)	0.2309 (3.937E-11)	0.2075 (3.826E-07)	-0.1693 (7.883E-15)	0.1164 (8.966E-22)	0.0024 (6.636E-12)
	10	0.1269 (2.634E-20)	0.1225 (5.655E-16)	0.0629 (1.191E-08)	-2.1277 (5.354E-25)	-5.3814 (5.946E-20)	-17.8760 (8.667E-12)
	20	0.3544 (1.580E-18)	0.3729 (2.536E-14)	0.3600 (1.228E-22)	-0.3647 (1.521E-26)	-1.4314 (8.279E-15)	-5.0739 (2.670E-13)
	all	0.3614 (6.777E-16)	0.3745 (2.961E-12)	0.3611 (1.637E-11)	-0.4204 (9.155E-26)	-1.4232 (1.170E-14)	-4.9260 (2.181E-13)
ANN	0	0.2772 (1.110E-16)	0.2693 (2.283E-14)	0.2415 (1.089E-08)	0.1924 (2.328E-25)	0.1193 (1.425E-19)	0.0177 (2.016E-10)
	10	0.3867 (9.693E-23)	0.3861 (2.699E-19)	0.3719 (1.481E-11)	-0.0020 (1.479E-22)	-0.3044 (1.383E-13)	-1.2751 (1.518E-11)
	20	0.2431 (7.053E-21)	0.2492 (1.389E-16)	0.2340 (5.689E-13)	-0.4173 (3.548E-04)	-0.9081 (6.624E-13)	-1.8306 (1.869E-11)
	all	-0.4727 (2.819E-15)	-0.5301 (5.803E-28)	-0.6294 (5.788E-20)	-1.0350 (8.204E-22)	-2.1418 (2.608E-13)	-5.1300 (3.830E-11)
RNN	0	0.3254 (8.035E-14)	0.3200 (1.731E-10)	0.3020 (8.619E-06)	0.1655 (1.473E-24)	0.0944 (8.473E-20)	0.0168 (2.083E-11)
	10	<b>0.4191</b> (2.713E-14)	<b>0.4254</b> (2.070E-11)	<b>0.4150</b> (1.021E-06)	0.3174 (2.404E-25)	0.0222 (1.345E-18)	-0.3291 (8.663E-12)
	20	0.3584 (5.590E-14)	0.3606 (5.548E-11)	0.3450 (1.817E-05)	0.3407 (1.316E-23)	0.1835 (3.676E-18)	-0.0238 (1.233E-12)
	all	0.3210 (2.498E-15)	0.3265 (8.128E-13)	0.2966 (2.166E-07)	<b>0.3722</b> (4.872E-25)	<b>0.1942</b> (2.693E-18)	-0.0373 (2.482E-10)

**Table 6. Weekly Forecasting Performance.**  $R^2_{OOS}$  scores for entire period and subperiod ( $t \in \{10, 5, 3, 2, 1\}$ ) are represented. The numbers in parentheses indicate the  $p$ -values of residual stationarity. The bold numbers are the best performance measurements for each subperiod.

**Forecasting Horizon** The longer the forecast horizon, the better the performances of most models tend to be (Table 4, 5, 6). This might be simply from the noise of the data, as data with shorter frequencies contain a smaller set of data for each realized volatility. However, on the contrary, the performance of the benchmark model decreases as the forecast horizon increases. Since the HAR model works well enough with the different horizons, it cannot be said that this is because memory is lost as the time horizon lengthens. Rather, it can be interpreted that the short-term variable and long-term variables of the HAR model do not sufficiently reflect the momentum of the realized volatility. We can see that the RFR and RNN models tend to make robust predictions over the forecast horizon compared to other models, and the results themselves are superior to others.

In addition to the forecast horizon, robustness to the sub-period length also has important implications. While it is very important to get it right over the entire time, it is also very important to get it right over the recent period. LASSO and ElasticNet show great performance for longer sub-periods, while RFR and RNN show better performance for recent times. The RFR and RNN show a very stable  $R^2$  score even during the period of COVID-19. This is particularly meaningful because not only the result is statistically significant (the residual is stationary) but they are robust to the test period. Interestingly, the results for LASSO and ElasticNet using recent time periods are a bit unusual. While the performance of the models themselves decreases as the forecast horizon decreases, the ranking of LASSO and ElasticNet actually increases. This suggests that oil prices have been very volatile in recent periods, and the momentum played much important role than other features.

Forecasting Model	Number of Feature Selection	$R^2_{OOS}$	$R^2_{OOS,10}$	$R^2_{OOS,5}$	$R^2_{OOS,3}$	$R^2_{OOS,2}$	$R^2_{OOS,1}$
HAR	-	-0.8224 (2.273E-03)	-0.8918 (4.362E-03)	-1.0895 (2.453E-01)	-0.1665 (9.286E-01)	-0.1312 (1.728E-02)	-1.0531 (7.422E-02)
TV-HAR	-	-0.1154 (5.893E-21)	-0.1759 (4.605E-19)	-0.2851 (1.943E-05)	-0.1238 (4.493E-01)	0.3608 (2.088E-01)	0.1883 (3.401E-02)
HAR-X	0	-0.8224 (2.273E-03)	-0.8918 (4.362E-03)	-1.0895 (2.453E-01)	-0.1665 (9.286E-01)	-0.1312 (1.728E-02)	-1.0531 (7.422E-02)
	10	0.1144 (3.581E-17)	0.0777 (1.620E-15)	0.0183 (1.738E-08)	-0.9011 (1.565E-02)	-1.0257 (3.521E-04)	-0.2328 (1.051E-02)
	20	0.0387 (1.409E-20)	0.0034 (1.080E-18)	-0.0793 (1.168E-10)	-1.3337 (4.948E-01)	-2.2149 (4.345E-06)	-0.3676 (4.611E-04)
	all	0.4685 (3.248E-07)	0.4636 (1.159E-13)	0.4581 (1.734E-04)	-1.1800 (2.572E-02)	-1.6785 (2.566E-01)	-0.4692 (1.000E+00)
		-0.0256 (5.237E-11)	-0.0788 (6.249E-10)	-0.1795 (4.716E-05)	-0.0518 (5.148E-01)	0.3053 (8.338E-01)	-0.0484 (7.259E-02)
LASSO	10	0.2071 (9.274E-13)	0.1678 (2.343E-11)	0.1058 (7.440E-02)	-0.0758 (1.269E-10)	0.3622 (1.406E-04)	0.2709 (3.048E-02)
	20	0.2215 (5.528E-13)	0.1830 (1.470E-11)	0.1116 (4.797E-02)	-0.0758 (1.269E-10)	0.3622 (1.406E-04)	0.2709 (3.048E-02)
	all	<b>0.5602</b> (1.104E-17)	<b>0.5474</b> (1.177E-15)	<b>0.5216</b> (2.227E-08)	-0.5571 (9.731E-02)	-0.2091 (5.403E-05)	0.0758 (7.254E-02)
		-0.0376 (1.311E-10)	-0.0911 (1.308E-09)	-0.1924 (8.332E-05)	-0.0595 (5.239E-01)	0.3015 (7.996E-01)	-0.0452 (6.511E-02)
	10	0.2301 (1.160E-12)	0.1929 (1.685E-11)	0.1344 (4.118E-02)	0.0662 (1.267E-01)	0.4722 (5.872E-05)	0.3081 (1.638E-02)
ElasticNet	20	0.2119 (8.128E-13)	0.1740 (1.224E-11)	0.1011 (4.506E-02)	0.0111 (1.356E-01)	0.4727 (5.603E-05)	0.3075 (1.653E-02)
	all	0.5581 (2.712E-17)	0.5463 (2.577E-15)	0.5212 (3.835E-08)	-0.5747 (8.831E-02)	-0.2567 (8.186E-05)	0.0769 (8.503E-02)
	0	-0.1846 (7.083E-12)	-0.2439 (1.508E-10)	-0.3223 (1.143E-05)	-2.5200 (6.053E-06)	-4.5855 (9.410E-01)	-9.4199 (4.531E-02)
	10	0.0772 (1.938E-13)	0.0730 (7.730E-12)	0.0865 (9.661E-02)	0.2048 (1.514E-06)	0.0403 (1.000E+00)	-0.1204 (3.485E-02)
	20	0.0855 (1.015E-13)	0.0826 (3.509E-12)	0.0670 (1.018E-01)	0.0598 (2.612E-09)	0.2555 (5.921E-01)	-0.2466 (9.447E-01)
DTR	all	-0.0329 (4.458E-13)	-0.0420 (1.372E-11)	-0.0856 (1.110E-01)	-0.0870 (1.333E-12)	0.1783 (2.495E-06)	-0.2413 (9.927E-01)
	0	0.1619 (7.406E-17)	0.1191 (3.315E-15)	0.0741 (7.310E-04)	-0.5734 (8.765E-09)	-0.2769 (8.249E-03)	-1.2444 (9.927E-01)
	10	0.3996 (8.176E-13)	0.3747 (1.728E-11)	0.3445 (9.641E-05)	0.2765 (5.663E-07)	<b>0.5252</b> (4.202E-04)	0.4003 (1.305E-01)
	20	0.3890 (2.925E-13)	0.3628 (5.877E-12)	0.3281 (7.337E-05)	0.2058 (1.140E-07)	0.4313 (1.715E-04)	<b>0.4681</b> (1.893E-01)
	all	0.3900 (3.774E-13)	0.3633 (7.489E-12)	0.3263 (1.052E-04)	<b>0.2798</b> (2.022E-01)	0.4507 (7.107E-01)	0.4183 (2.762E-01)
GBR	0	0.1129 (1.800E-17)	0.0680 (8.916E-16)	0.0209 (2.536E-04)	-0.5325 (8.862E-08)	-0.3486 (4.154E-02)	-1.1792 (1.000E+00)
	10	0.3156 (6.920E-13)	0.2944 (1.740E-11)	0.2640 (6.580E-05)	0.1661 (3.576E-01)	0.2838 (1.597E-04)	0.2808 (2.158E-01)
	20	0.3279 (1.122E-11)	0.3078 (1.901E-10)	0.2804 (2.462E-05)	-0.6926 (7.473E-08)	-1.2274 (2.856E-09)	0.3230 (9.529E-02)
	all	0.3690 (2.175E-12)	0.3447 (4.404E-11)	0.3150 (9.949E-06)	-0.4082 (1.114E-07)	-0.3846 (2.762E-07)	0.3471 (1.276E-01)
		-0.1702 (7.448E-12)	-0.2119 (1.009E-11)	-0.3231 (8.182E-05)	-0.0139 (6.271E-01)	0.2969 (7.602E-01)	0.0773 (1.347E-01)
ANN	10	0.4105 (6.680E-14)	0.3971 (4.604E-13)	0.3884 (5.001E-05)	-0.1129 (2.673E-01)	0.1507 (1.000E+00)	-0.1438 (6.024E-01)
	20	0.3738 (1.097E-16)	0.3486 (1.992E-15)	0.3458 (1.470E-16)	-0.1116 (2.468E-07)	-0.2108 (7.774E-01)	-0.6445 (6.415E-01)
	all	0.3484 (5.600E-21)	0.3289 (1.795E-19)	0.3275 (6.915E-04)	-0.7419 (2.500E-07)	-0.9420 (4.844E-01)	-2.6058 (7.730E-01)
	0	0.0469 (3.390E-10)	0.0066 (2.627E-09)	-0.0369 (7.359E-04)	-0.0439 (6.892E-01)	0.1248 (4.299E-01)	-0.4881 (3.271E-01)
	10	0.1347 (4.888E-10)	0.0939 (6.055E-09)	0.0426 (7.142E-04)	0.1647 (3.782E-08)	0.4339 (1.713E-02)	0.2462 (1.278E-01)
RNN	20	0.1524 (1.322E-10)	0.1202 (1.146E-09)	0.0652 (1.569E-03)	0.0964 (2.063E-01)	0.1485 (1.321E-02)	-0.0698 (9.962E-01)
	all	0.0755 (7.702E-10)	0.0408 (2.657E-09)	0.0187 (2.060E-04)	0.2765 (2.970E-07)	0.5025 (1.557E-02)	0.4409 (9.990E-01)
	0	0.0469 (3.390E-10)	0.0066 (2.627E-09)	-0.0369 (7.359E-04)	-0.0439 (6.892E-01)	0.1248 (4.299E-01)	-0.4881 (3.271E-01)
	10	0.1347 (4.888E-10)	0.0939 (6.055E-09)	0.0426 (7.142E-04)	0.1647 (3.782E-08)	0.4339 (1.713E-02)	0.2462 (1.278E-01)
	20	0.1524 (1.322E-10)	0.1202 (1.146E-09)	0.0652 (1.569E-03)	0.0964 (2.063E-01)	0.1485 (1.321E-02)	-0.0698 (9.962E-01)

**Table 7. Monthly Forecasting Performance with Uncertainty Indices.**  $R^2_{OOS}$  scores for entire period and subperiod ( $t \in \{10, 5, 3, 2, 1\}$ ) are represented. The numbers in parentheses indicate the  $p$ -values of residual stationarity. The bold numbers are the best performance measurements for each subperiod.

Forecasting Model	Number of Feature Selection	$R^2_{OOS}$	$R^2_{OOS,10}$	$R^2_{OOS,5}$	$R^2_{OOS,3}$	$R^2_{OOS,2}$	$R^2_{OOS,1}$
HAR	-	-0.4331 (1.471E-15)	-0.5165 (2.663E-12)	-0.6729 (2.883E-07)	0.0788 (3.739E-12)	0.0702 (2.547E-03)	-0.0623 (2.408E-06)
TV-HAR	-	0.1240 (8.207E-18)	0.0891 (7.085E-14)	0.0197 (1.216E-07)	0.2160 (1.090E-12)	0.2571 (1.267E-12)	0.1236 (6.754E-06)
HAR-X	0	-0.4331 (1.471E-15)	-0.5165 (2.663E-12)	-0.6729 (2.883E-07)	0.0788 (3.739E-12)	0.0702 (2.547E-03)	-0.0623 (2.408E-06)
	10	0.2601 (4.170E-14)	0.2612 (9.497E-12)	0.2202 (2.150E-06)	-0.2401 (3.988E-13)	-0.0644 (4.247E-13)	-0.1814 (5.232E-01)
	20	0.2335 (1.692E-13)	0.2304 (4.809E-11)	0.1916 (3.602E-06)	-0.5422 (1.840E-12)	-0.1445 (7.702E-13)	0.0414 (3.656E-07)
	all	0.2516 (7.572E-10)	0.2458 (4.139E-08)	0.1988 (6.211E-05)	-0.7750 (1.112E-11)	-0.9573 (3.391E-09)	-0.5857 (5.791E-04)
LASSO	0	0.1724 (3.072E-17)	0.1455 (4.827E-14)	0.0814 (1.725E-08)	0.2501 (1.690E-12)	0.2458 (8.294E-12)	0.0157 (1.583E-05)
	10	0.2412 (4.763E-17)	0.2188 (4.405E-14)	0.1670 (2.112E-08)	0.1928 (1.998E-11)	0.2668 (3.151E-12)	0.0814 (3.661E-05)
	20	0.2409 (4.693E-17)	0.2188 (4.405E-14)	0.1670 (2.112E-08)	0.1928 (1.998E-11)	0.2668 (3.151E-12)	0.0814 (3.661E-05)
	all	0.3177 (3.398E-12)	0.3053 (1.112E-09)	0.2517 (2.688E-05)	0.1807 (2.919E-14)	-0.0947 (3.681E-12)	0.0516 (7.178E-05)
ElasticNet	0	0.1713 (2.450E-17)	0.1451 (3.507E-14)	0.0810 (1.360E-08)	0.2486 (1.721E-12)	0.2446 (8.137E-12)	0.0197 (1.713E-05)
	10	0.2844 (1.273E-16)	0.2651 (3.364E-14)	0.2191 (2.301E-08)	0.2263 (1.449E-11)	0.2989 (5.154E-13)	0.2097 (2.380E-05)
	20	0.2814 (7.886E-17)	0.2602 (3.258E-14)	0.2126 (2.291E-08)	0.2216 (1.387E-11)	0.2985 (5.136E-13)	0.2097 (2.380E-05)
	all	0.3094 (1.226E-12)	0.2958 (5.212E-10)	0.2429 (1.578E-05)	0.0689 (1.446E-13)	-0.2175 (2.824E-11)	-0.0307 (1.611E-04)
DTR	0	-0.4030 (2.277E-15)	-0.4207 (7.154E-12)	-0.4609 (1.433E-07)	-1.0836 (3.983E-14)	-1.4369 (3.434E-12)	-2.4828 (9.017E-06)
	10	-0.0261 (4.796E-17)	-0.0170 (5.623E-13)	-0.0362 (7.743E-07)	-1.8843 (1.928E-16)	-0.7706 (9.395E-17)	-1.0411 (6.755E-03)
	20	0.1224 (3.113E-11)	0.1261 (4.435E-09)	0.1072 (3.469E-05)	-1.8057 (5.907E-15)	-0.4107 (1.873E-12)	-0.7176 (9.336E-05)
	all	0.0022 (7.695E-11)	-0.0103 (4.156E-08)	-0.0634 (1.761E-04)	-0.0145 (1.140E-07)	-0.1713 (4.450E-16)	-0.0195 (1.961E-09)
RFR	0	0.2228 (8.574E-15)	0.2229 (5.530E-12)	0.1903 (1.257E-06)	-0.3520 (1.421E-15)	0.0430 (1.502E-12)	0.0110 (1.322E-01)
	10	0.4157 (5.086E-13)	0.4263 (1.075E-10)	0.3989 (2.047E-06)	0.4385 (1.168E-12)	0.2313 (1.199E-08)	0.0671 (3.233E-05)
	20	0.4332 (3.607E-15)	0.4365 (1.372E-12)	0.4103 (1.168E-07)	0.4379 (3.055E-12)	0.2529 (3.044E-12)	0.1195 (8.426E-06)
	all	0.4220 (9.403E-16)	0.4262 (3.742E-13)	0.3991 (6.208E-08)	0.4003 (1.501E-12)	0.2742 (2.426E-12)	0.2264 (9.950E-06)
GBR	0	0.1325 (4.129E-15)	0.1324 (3.481E-12)	0.0997 (6.767E-07)	-0.6054 (1.188E-15)	0.0271 (1.539E-12)	-0.1844 (9.086E-02)
	10	0.4224 (4.113E-12)	0.4355 (2.471E-10)	0.4149 (2.917E-06)	0.4191 (3.185E-10)	0.3391 (3.119E-12)	0.1340 (2.318E-05)
	20	<b>0.4382</b> (1.572E-14)	<b>0.4516</b> (2.465E-12)	<b>0.4395</b> (8.803E-08)	0.3489 (4.628E-08)	0.2681 (1.491E-10)	0.1326 (3.073E-06)
	all	0.4275 (4.274E-15)	0.4409 (4.644E-13)	0.4327 (2.171E-08)	0.4118 (3.122E-11)	0.2815 (2.002E-11)	0.2174 (2.143E-05)
ANN	0	0.0682 (1.534E-17)	0.0427 (4.212E-14)	-0.0317 (4.223E-09)	0.1667 (7.337E-14)	0.0431 (4.340E-13)	-0.1548 (1.125E-05)
	10	0.2614 (2.180E-10)	0.2473 (1.381E-09)	0.1919 (1.600E-05)	0.1750 (8.319E-09)	0.1655 (8.225E-11)	-0.0554 (1.199E-03)
	20	0.3023 (2.594E-12)	0.2989 (1.920E-11)	0.2736 (2.379E-06)	0.0370 (6.188E-12)	-0.1053 (7.267E-11)	-0.1325 (9.531E-06)
	all	0.2774 (7.206E-14)	0.2871 (2.639E-12)	0.2747 (1.778E-07)	-0.3376 (2.387E-10)	-0.4558 (1.255E-03)	-0.3589 (1.872E-01)
RNN	0	0.2079 (1.431E-14)	0.1921 (8.350E-12)	0.1485 (1.880E-06)	0.1730 (6.609E-14)	0.0915 (1.368E-12)	-0.0419 (2.223E-06)
	10	0.3486 (4.026E-13)	0.3379 (1.103E-10)	0.3027 (2.899E-06)	0.4052 (9.886E-09)	<b>0.3539</b> (4.915E-11)	0.1496 (3.805E-03)
	20	0.3111 (7.085E-14)	0.3055 (1.848E-11)	0.2787 (8.996E-07)	0.3561 (1.002E-07)	0.2412 (1.206E-01)	0.2320 (1.697E-04)
	all	0.2516 (3.833E-13)	0.2280 (8.594E-11)	0.1828 (2.698E-06)	<b>0.4650</b> (3.253E-09)	0.2926 (2.174E-03)	<b>0.2460</b> (2.352E-04)

**Table 8. Bi-weekly Forecasting Performance with Uncertainty Indices.**  $R^2_{OOS}$  scores for entire period and subperiod ( $t \in \{10, 5, 3, 2, 1\}$ ) are represented. The numbers in parentheses indicate the  $p$ -values of residual stationarity. The bold numbers are the best performance measurements for each subperiod.



Forecasting Model	Number of Feature Selection	$R^2_{OOS}$	$R^2_{OOS,10}$	$R^2_{OOS,5}$	$R^2_{OOS,3}$	$R^2_{OOS,2}$	$R^2_{OOS,1}$
HAR	-	-0.0714 (6.045E-17)	-0.1217 (9.597E-13)	-0.2127 (2.885E-07)	-0.0059 (6.728E-18)	0.0323 (7.068E-17)	-0.1408 (5.059E-09)
TV-HAR	-	0.3117 (2.151E-19)	0.2990 (6.913E-15)	0.2749 (8.617E-09)	0.1455 (6.334E-27)	0.1225 (8.163E-11)	0.0129 (1.361E-14)
HAR-X	0	-0.0714 (6.045E-17)	-0.1217 (9.597E-13)	-0.2127 (2.885E-07)	-0.0059 (6.728E-18)	0.0323 (7.068E-17)	-0.1408 (5.059E-09)
	10	0.3059 (1.592E-16)	0.3186 (5.208E-14)	0.2995 (4.284E-08)	-0.0861 (1.004E-23)	0.0106 (2.169E-20)	-0.1231 (1.450E-12)
	20	0.2824 (7.778E-16)	0.3074 (4.292E-13)	0.2823 (1.648E-07)	-0.2541 (8.955E-25)	-0.0957 (7.375E-22)	0.0305 (7.469E-14)
	all	0.3474 (4.481E-11)	0.3852 (4.360E-08)	0.3628 (8.877E-05)	-0.2546 (4.451E-23)	-0.2978 (2.499E-17)	-0.2114 (7.067E-12)
LASSO	0	0.3367 (1.174E-14)	0.3321 (1.205E-11)	0.3114 (8.683E-07)	0.1872 (1.618E-26)	0.1673 (2.353E-22)	0.0231 (1.272E-14)
	10	0.3376 (1.282E-14)	0.3331 (1.221E-11)	0.3126 (8.546E-07)	0.1889 (2.387E-26)	0.1700 (3.201E-22)	0.0260 (1.472E-14)
	20	0.3376 (1.283E-14)	0.3331 (1.221E-11)	0.3126 (8.546E-07)	0.1889 (2.387E-26)	0.1700 (3.201E-22)	0.0260 (1.472E-14)
	all	0.3848 (9.367E-14)	0.3969 (7.932E-11)	0.3705 (3.422E-06)	0.1992 (7.732E-25)	0.0852 (2.885E-05)	0.0585 (9.827E-14)
ElasticNet	0	0.3398 (1.618E-16)	0.3332 (5.191E-13)	0.3124 (1.502E-07)	0.1852 (1.320E-26)	0.1645 (1.758E-22)	0.0248 (1.237E-14)
	10	0.3588 (7.259E-16)	0.3524 (1.281E-12)	0.3353 (1.967E-07)	0.2128 (4.884E-26)	0.2001 (5.583E-22)	0.0822 (3.547E-14)
	20	0.3575 (6.743E-16)	0.3513 (1.131E-12)	0.3340 (1.810E-07)	0.1981 (3.484E-26)	0.1948 (1.171E-22)	0.0977 (1.123E-14)
	all	0.3878 (9.756E-14)	0.3984 (1.049E-10)	0.3735 (3.877E-06)	0.1769 (3.044E-25)	0.0863 (1.489E-05)	0.0564 (3.906E-14)
DTR	0	-0.2966 (2.571E-18)	-0.2424 (9.006E-14)	-0.2173 (7.943E-08)	-0.7698 (1.270E-21)	-1.2630 (1.293E-10)	-1.5798 (1.565E-10)
	10	-0.2386 (4.435E-18)	-0.2220 (1.048E-13)	-0.2806 (8.784E-08)	-0.1886 (6.482E-25)	-0.7208 (8.301E-17)	-1.1564 (5.378E-12)
	20	-0.2525 (2.395E-12)	-0.2309 (6.108E-09)	-0.2726 (1.036E-04)	-0.3411 (3.523E-29)	-0.4462 (1.706E-10)	-0.7275 (1.472E-12)
	all	-0.3255 (4.783E-14)	-0.3111 (1.652E-08)	-0.3645 (2.058E-08)	-0.5204 (1.362E-29)	-0.7524 (1.248E-03)	-0.9286 (3.870E-08)
RFR	0	0.2455 (4.173E-16)	0.2510 (1.282E-12)	0.2449 (2.246E-07)	0.0316 (5.965E-25)	0.0896 (4.826E-10)	-0.0775 (5.790E-09)
	10	0.4214 (2.654E-15)	<b>0.4407</b> (1.983E-12)	<b>0.4274</b> (1.080E-07)	0.3240 (4.968E-24)	0.1997 (1.325E-18)	0.0566 (1.695E-10)
	20	<b>0.4227</b> (1.143E-16)	0.4376 (2.685E-13)	0.4243 (4.200E-08)	0.3481 (7.566E-26)	0.1747 (1.782E-05)	<b>0.1215</b> (1.674E-12)
	all	0.4206 (6.574E-17)	0.4353 (1.226E-13)	0.4196 (1.633E-08)	0.3022 (5.679E-26)	0.1179 (2.372E-05)	0.0993 (1.249E-12)
GBR	0	0.2001 (2.660E-14)	0.2083 (1.094E-12)	0.1802 (1.146E-06)	-0.1197 (1.383E-14)	0.1183 (1.108E-21)	0.0202 (1.710E-11)
	10	0.3083 (1.613E-16)	0.3169 (2.175E-11)	0.2931 (1.079E-06)	0.2776 (8.415E-24)	0.1671 (8.784E-20)	0.0378 (1.140E-09)
	20	0.3561 (1.443E-16)	0.3694 (6.048E-12)	0.3557 (8.274E-07)	0.2418 (2.766E-23)	0.1512 (3.112E-05)	0.0228 (1.756E-11)
	all	0.3497 (1.997E-15)	0.3606 (3.194E-12)	0.3395 (5.048E-07)	0.2606 (2.245E-23)	0.1617 (2.052E-05)	0.0634 (4.427E-12)
ANN	0	0.2995 (2.356E-16)	0.2950 (2.791E-14)	0.2730 (1.293E-08)	0.2002 (1.167E-25)	0.1233 (4.722E-20)	0.0065 (2.788E-11)
	10	0.3883 (2.275E-15)	0.3957 (4.064E-12)	0.3864 (1.313E-06)	0.2196 (6.690E-18)	0.1475 (1.359E-16)	-0.1258 (2.137E-08)
	20	0.2902 (1.299E-16)	0.2980 (2.206E-13)	0.2911 (3.362E-07)	-0.1979 (9.375E-23)	-0.2956 (4.379E-21)	-1.0640 (2.091E-12)
	all	0.2713 (1.271E-15)	0.2888 (3.977E-12)	0.3015 (1.706E-06)	-0.1994 (1.129E-21)	-0.6120 (1.359E-03)	-1.4317 (3.721E-09)
RNN	0	0.3247 (6.578E-14)	0.3203 (1.723E-10)	0.3034 (8.985E-06)	0.2132 (1.041E-24)	0.1238 (5.747E-21)	0.0993 (1.853E-11)
	10	0.4174 (3.292E-13)	0.4222 (2.262E-10)	0.4135 (4.210E-06)	0.3297 (1.956E-23)	0.1765 (4.850E-19)	0.0408 (5.604E-09)
	20	0.3895 (3.013E-15)	0.4044 (2.065E-12)	0.3974 (1.171E-07)	0.2838 (1.568E-24)	0.0861 (1.835E-14)	-0.1597 (5.615E-10)
	all	0.3473 (5.470E-15)	0.3540 (1.796E-12)	0.3236 (4.216E-07)	<b>0.4095</b> (1.039E-23)	<b>0.2544</b> (1.597E-19)	-0.0095 (4.046E-11)

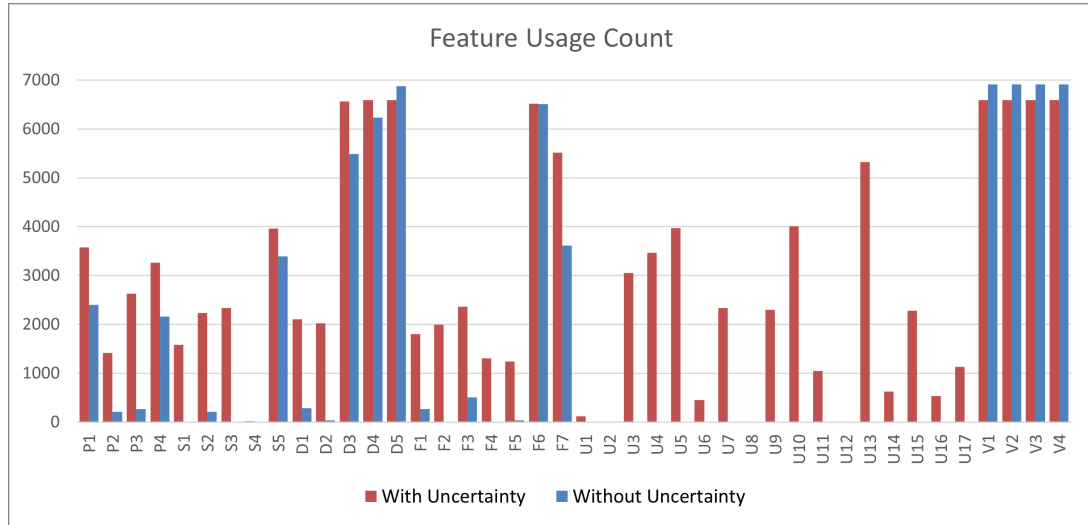
**Table 9. Weekly Forecasting Performance with Uncertainty Indices.**  $R^2_{OOS}$  scores for entire period and subperiod ( $t \in \{10, 5, 3, 2, 1\}$ ) are represented. The numbers in parentheses indicate the  $p$ -values of residual stationarity. The bold numbers are the best performance measurements for each subperiod.

**Uncertainty Indices** Furthermore, we find that the use of uncertainty indices can push the accuracy of models even higher (Table 7, 8, and 9). The results vary slightly by model, with ElasticNet, DTR, and RNN seeing performance improvements by adding uncertainty features in about half of the cases, while LASSO, RFR, and ANN improve in most cases.

There are also significant differences based on the forecasting horizon. For most models, the shorter the forecast horizon, the better the use of uncertainty indices. This can be explained by the fact that the uncertainty indices reduce overfitting and are more influential in dealing with noise in short-term forecasts.

In many cases, adding uncertainty data to the training data set is beneficial. However, we could not guarantee that all uncertainty data has high explanatory power because many of them are likely to be pruned during the feature selection process. Even so, the highest performer of each model shows a higher performance when uncertainty indices join the training data set.

### Feature Selection



**Figure 2. Feature Usage Count** The feature usage was aggregated for the best performance models for each time horizon that selected 10 or 20 features, excluding cases where the entire set of features was selected. Each feature is assigned a number in Table 1 based on the first letter of each factor group (P, S, D, F, U, V) and in order from top to bottom (except for WTI future price, which is used as a dependent variable), e.g., WTI Spot price is equal to P2, PPI in US is equal to D4, etc.

We examine the feature usage in the models. First, we defined the best performance model as the feature selection methodology that performed the best in each model for

each time horizon being predicted. Then, we counted and summarized the usage of the features utilized by each best model when predicting each time point. Also, in order to clearly see the difference in the number of features selected, we excluded all cases where all features were selected and only kept cases where 10 or 20 features were selected. The result is presented in Figure 2. Each feature is assigned a number in Table 1 based on the first letter of each factor group (P, S, D, F, U, V) and in order from top to bottom (except for WTI future price, which is used as a dependent variable). For example, WTI Spot price is equal to P2, PPI in US is equal to D4.

Looking at the Figure 2, we can see that several features are highly selected, both with and without uncertainty features. First of all, the volatility indices(V1 to 4) are all highly selected, with PPI in China, US, and EU(D3 to 5) and federal funds rate and MSCI World Standard Index(F6,7) being the most selected. Next, WTI spot prices and NGL spot prices(P1,4) and capacity utilization rate(S5) are being selected. For the uncertainty features, US economic policy uncertainty in financial regulation, monetary policy, national security, and economic policy uncertainty(U13,5,10,4) are mostly selected. As in Degiannakis and Filis (2022) and Delis et al. (2023), the implied volatility indices are important in forecasting the volatility. Also, as the oil is one of the most important commodity in production, the PPI plays significant role. Federal funds rate and MSCI World Index represents the global economy, which is important for the volatility.

On the other hand, total OPEC production capacity(S4) and Brent oil spot prices(P2), US/UK, China/UK foreign exchange rate(F4,5) failed to being selected. Also, Global economic policy uncertainty, US economic policy uncertainty in fiscal policy, government spending, regulation, trade policy, and World uncertainty index(U1,2,6,8,12,14,16) are not selected.

Interestingly, there were features that were selected in the models with uncertainty but not in the models without uncertainty. These features are all price features(P1 to 4), global crude oil production, stock, and export(S1 to 3), MSCI World Standard Index(F7), and global crude oil import and liquid fuels consumption(D1,2). These differences are expected to have contributed to the difference in performance between the two models.

**Potential Explanation** Now, two main question rises. Why the performance orders are different for each forecast horizon? How do some machine learning models outperform the linear model? To answer the first question, we first see the performance of lagged value, in other words, the No-change model. The no-change model calculates the  $R^2$  score of only the first lagged value of the target time series. Therefore, it calculates the effect of autocorrelation of lag 1. For many cases, the No-change model's  $R^2$  is positive, which means the autocorrelation's effect is not negligible. Because the models with different forecast horizons are exposed to the different levels of autocorrelation effect, the performance orders could change.

The answer to the second question can come from many points, but the assumption of Gaussian distribution seems most restricted to the linear models. Our linear model assumes that the residual follows the normal distribution. However, because the model is linear, the real-world fat-tailed data may not fit the linear model. The nonlinear model, on the other hand, has the potential to explain such fat-tailed data. Another reason is the in-sample over-fitting. For example, the decision tree model is notorious for its tendency to over-fit. To overcome that issue, researchers develop a method to reduce the shortcoming and one way is the random forest, one of our best performers. Last but not least, the model may not converge well. Data is not always clean. Such ill-posed data can hurt the model during the training stage which finds the optimal weights to minimize the error.

#### 4.3. *Discussion and Possible Extensions*

Volatility is one of the key variables for trading strategies, asset allocation, and risk management. Now we introduce the quantitative analysis suggested by previous studies on the usefulness of volatility forecasting. One way to express the quantitative advantage for accurate future volatility is the certainty equivalent return (CER). Yin and Yang (2016), Ma et al. (2018), Zhang et al. (2018), and Zhang et al. (2019) show that the mean-variance utility investor can benefit from forecasting future volatility in terms of CER. When the mean-variance utility investor considers the crude oil future and

risk-free asset, the portfolio return is

$$R_p = wl(r + r_f) + (1 - w)r_f$$

where  $w$  is the portfolio weight of crude oil futures,  $l$  is the leverage ratio,  $r$  is the excess return, and  $r_f$  is the risk-free rate. Then, CER of the investor is:

$$CER = \bar{R}_p - \frac{1}{2}\gamma\sigma_p^2$$

where  $\bar{R}_p$  and  $\sigma_p^2$  are the mean and the variance of the portfolio return  $R_p$  over the out-of-sample period. It is straightforward that a more accurate prediction yields a higher CER.

Another way to emphasize the importance of a volatility forecast is the Sharpe ratio of the simple directional trading strategy. Following Easley et al. (2020), consider  $n$  independent and identically distributed bets, where the outcome  $y$  of a bet with a profit  $\pi > 0$  with probability  $P[y = \pi] = p$  and a loss  $\pi$  with probability  $1 - p$ . The expected profit is  $E[y] = \pi(2p - 1)$ , and the variance is  $Var[y] = 4\pi^2 p(1 - p)$ . Then, the Sharpe ratio is:

$$\theta(p, n) = \frac{nE[y]}{\sqrt{nVar[y]}} = \frac{2p - 1}{2\sqrt{p(1 - p)}}\sqrt{n}$$

Easley et al. (2020) stated that repeated trading with sufficiently large enough trials ( $n = 13,000$ ) can achieve the Sharpe ratio 2.04 with the small enhancement of directional forecast  $p = 0.52$ . However, our models achieve far more accurate predictability around 0.6 to 0.7. Therefore, even smaller trials  $n = 100$  can get a higher Sharpe ratio of 4.36 with  $p = 0.7$ .

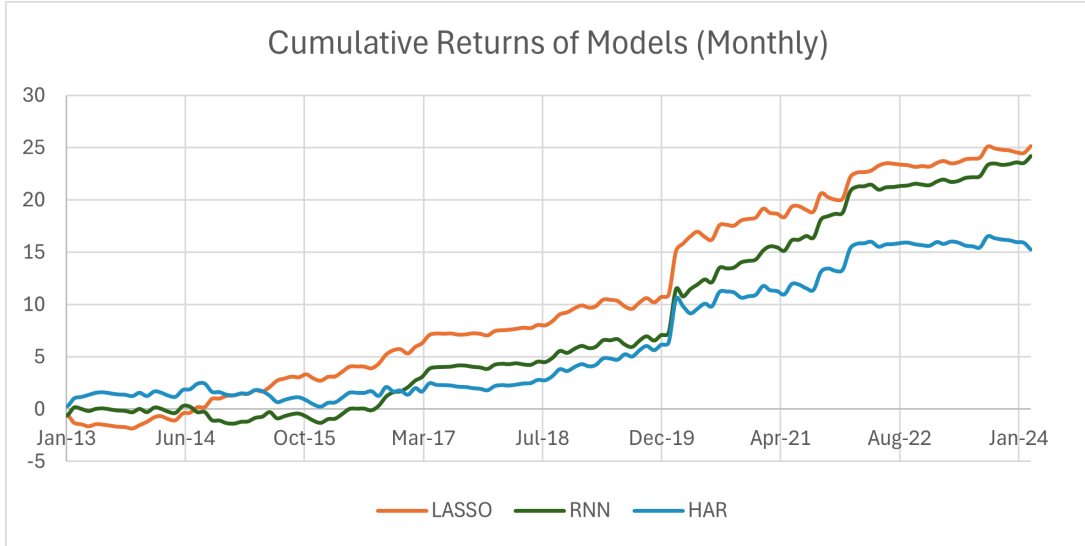
The bet above is to analyze the advantage of accurate directional forecasting. However, in practice, we can consider the investment in the volatility ETFs such as OVX or trading strategies based on volatility. Delis et al. (2022) introduces a simple strategy for volatility trading. In the case where the  $s$ -step forecasted oil price volatility of the model  $i$  at time  $t$  (denoted as  $\widehat{RV}_{t+s|t}$ ) is higher than that of the actual volatility at time

$t$  (denoted as  $RV_t$ ), the trader will take a long position in an asset that exhibits similar performance to the oil realized volatility, such as OVX. Conversely, if the forecasted volatility of model  $i$  at time  $t + s$  is lower than that of the actual volatility at time  $t$ , the trader will take a short position. Consequently, the cumulative returns of model  $i$ ,  $r^i$  over the out-of-sample forecasting period 1 to  $t_N$  can be calculated as follows:

$$r^i = \sum_{t=1}^{t_N-s} \frac{(RV_{t+s} - RV_t)d_t^i}{RV_t}$$

where  $d_t^i = 1$  if  $\widehat{RV}_{t+s|t} > RV_t$ , and  $d_t^i = -1$  if  $\widehat{RV}_{t+s|t} \leq RV_t$ .

Over the past 11 years, HAR has generated cumulative returns of 1,589%, while LASSO and RNN have generated cumulative returns of 2,448% and 2,355% over the same period. It is clear that the machine learning models outperform the benchmark HAR, with LASSO demonstrating the most impressive performance overall. However, for RNN, which exhibited superior results in the recent period, the cumulative return increase exhibits a steeper slope than LASSO, and this period also encompasses the COVID-19 global pandemic, suggesting that RNN models are also significant. We showed the cumulative returns of LASSO, RNN, and benchmark HAR in Figure 3 below.



**Figure 3. Cumulative Returns of LASSO, RNN, and HAR** The forecasting horizon is monthly, and both models contain uncertainty factors. Furthermore, both models utilize all features selected through feature selection.

The preceding examples illustrate the effectiveness of a precise volatility forecast

in various contexts. This paper demonstrates that integrating representative machine learning models with the conventional HAR model can significantly enhance forecasting performance.

Going further, we might be able to improve performance through various detailed methodologies for time-series forecasting. There can be a variety of improvements - for instance, combining multiple forecasting models such as dynamic model averaging, dynamic model selection, or other ensemble methods (Wang et al. (2016); Audrino and Knaus (2016); Ding (2018); Zhang et al. (2018)), improving feature selection methods (Ghadimi et al. (2018); Karasu et al. (2020)), and reducing the noise of time series data using bilateral filters, wavelet denoising, autoencoder, etc (Uddin et al. (2019); Wang et al. (2017)).

## 5. Conclusion

This paper compares and analyzes the predictability of the realized volatility of crude oil future prices with various forecasting models from April 2002 to April 2024 including the Great Recession and COVID-19. With a wide variety of analyses and comparisons, we can explore the possibility of the potential usage of machine learning models in the field of volatility forecasting. A large set of explanatory variables is considered, and each explanatory variable belongs to one of six groups -prices, supply, demand, financial, implied volatility indices, and uncertainty factors. The walk-forward cross-validation shows the out-of-sample forecasting performance of ten types of forecasting models - HAR, TV-HAR, HAR-X, LASSO, ElasticNet, DTR, RFR, GBR, ANN, and RNN.

Although HAR is a popular conventional model in forecasting the realized volatility, the performance decreases in out-of-sample. Various machine learning models with the momentum factors of the HAR were tested, and some of them have significantly outperformed out-of-sample forecasts, such as the RFR and RNN.

Additionally, an investigation was conducted into the features selected in the models, which revealed that the implied volatility indices and uncertainty factors were the most frequently chosen. Furthermore, the PPI in China, the US, and the EU, the federal funds rate, and the MSCI World Standard Index were identified as the most frequently

selected variables.

Diverse previous literature has dealt with the advantages of accurate volatility forecasting. Volatility is one of the key variables for trading strategies, asset allocation, and risk management. Therefore, future works to enhance the predictive power are worthy with numerous potential extensions and applications.

### **Disclosure statement**

There is no relevant interest to declare.

### **Funding**

This work was supported by the Ministry of Education of the Republic of Korea and the National Research Foundation of Korea.(NRF-2022S1A3A2A02089950)



## References

- Ahir, H., N. Bloom, and D. Furceri (2018). The world uncertainty index. *Working Paper*, Available at SSRN 3275033.
- Apergis, N. and S. M. Miller (2009). Do structural oil-market shocks affect stock prices? *Energy economics* 31(4), 569–575.
- Audrino, F. and S. D. Knaus (2016). Lassoing the har model: A model selection perspective on realized volatility dynamics. *Econometric Reviews* 35(8-10), 1485–1521.
- Baker, S. R., N. Bloom, and S. J. Davis (2016). Measuring economic policy uncertainty. *The Quarterly Journal of Economics* 131(4), 1593–1636.
- Baker, S. R., N. Bloom, S. J. Davis, and K. J. Kost (2019). Policy news and stock market volatility. Technical report, National Bureau of Economic Research.
- Baker, S. R., N. Bloom, S. J. Davis, K. J. Kost, M. C. Sammon, and T. Viratyosin (2020). The unprecedented stock market impact of covid-19. Technical report, National Bureau of Economic Research.
- Brown, S. P. and M. K. Yucel (2008). What drives natural gas prices? *The Energy Journal* 29(2).
- Caggiano, G., E. Castelnuevo, and R. Kima (2020). The global effects of covid-19-induced uncertainty. *Economics Letters* 194(109392).
- Christiansen, C., M. Schmeling, and A. Schrimpf (2012). A comprehensive look at financial volatility prediction by economic variables. *Journal of Applied Econometrics* 27(6), 956–977.
- Corsi, F. (2009). A simple approximate long-memory model of realized volatility. *Journal of Financial Econometrics* 7(2), 174–196.
- Davis, S. J. (2016). An index of global economic policy uncertainty. Technical report, National Bureau of Economic Research.
- Degiannakis, S. and G. Filis (2022). Oil price volatility forecasts: What do investors need to know? *Journal of International Money and Finance* 123, 102594.
- Delis, P., S. Degiannakis, and G. Filis (2022). What matters when developing oil price volatility forecasting frameworks? *Journal of Forecasting* 41(2), 361–382.
- Delis, P., S. Degiannakis, and K. Giannopoulos (2023). What should be taken into consideration when forecasting oil implied volatility index? *The Energy Journal* 44(5), 231–250.
- Diebold, F. X. and R. S. Mariano (2002). Comparing predictive accuracy. *Journal of Business & Economic Statistics* 20(1), 134–144.
- Ding, Y. (2018). A novel decompose-ensemble methodology with aic-ann approach for crude

- oil forecasting. *Energy* 154, 328–336.
- Easley, D., M. López de Prado, M. O’Hara, and Z. Zhang (2020, 07). Microstructure in the Machine Age. *The Review of Financial Studies* 34(7), 3316–3363.
- Ghadimi, N., A. Akbarimajd, H. Shayeghi, and O. Abedinia (2018). Two stage forecast engine with feature selection technique and improved meta-heuristic algorithm for electricity load forecasting. *Energy* 161, 130–142.
- Ghoddusi, H., G. G. Creamer, and N. Rafizadeh (2019). Machine learning in energy economics and finance: A review. *Energy Economics* 81, 709–727.
- Goodfellow, I. (2016). *Deep learning*, Volume 196. MIT press.
- Gu, S., B. Kelly, and D. Xiu (2020). Empirical asset pricing via machine learning. *The Review of Financial Studies* 33(5), 2223–2273.
- Guo, B. and H. Lin (2020). Volatility and jump risk in option returns. *Journal of Futures Markets* 40(11), 1767–1792.
- Hallock Jr, J. L., P. J. Tharakan, C. A. Hall, M. Jefferson, and W. Wu (2004). Forecasting the limits to the availability and diversity of global conventional oil supply. *Energy* 29(11), 1673–1696.
- Hamilton, J. D. (2009). Causes and consequences of the oil shock of 2007-08. Technical report, National Bureau of Economic Research.
- Hammoudeh, S. M., B. T. Ewing, and M. A. Thompson (2008). Threshold cointegration analysis of crude oil benchmarks. *The Energy Journal* 29(4).
- Hansen, P. R., A. Lunde, and J. M. Nason (2011). The model confidence set. *Econometrica* 79(2), 453–497.
- Harvey, D., S. Leybourne, and P. Newbold (1997). Testing the equality of prediction mean squared errors. *International Journal of forecasting* 13(2), 281–291.
- Hastie, T., R. Tibshirani, and J. Friedman (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media.
- Herrera, A. M., M. B. Karaki, and S. K. Rangaraju (2019). Oil price shocks and us economic activity. *Energy Policy* 129, 89–99.
- Karasu, S., A. Altan, S. Bekiros, and W. Ahmad (2020). A new forecasting model with wrapper-based feature selection approach using multi-objective optimization technique for chaotic crude oil time series. *Energy* 212, 118750.
- Kilian, L. (2009). Not all oil price shocks are alike: Disentangling demand and supply shocks in the crude oil market. *American Economic Review* 99(3), 1053–69.

- Kilian, L. and B. Hicks (2013). Did unexpectedly strong economic growth cause the oil price shock of 2003–2008? *Journal of Forecasting* 32(5), 385–394.
- Ma, F., J. Liu, M. Wahab, and Y. Zhang (2018). Forecasting the aggregate oil price volatility in a data-rich environment. *Economic Modelling* 72, 320–332.
- Ma, F., M. Wahab, D. Huang, and W. Xu (2017). Forecasting the realized volatility of the oil futures market: A regime switching approach. *Energy Economics* 67, 136–145.
- Miao, H., S. Ramchander, T. Wang, and D. Yang (2017). Influential factors in crude oil price forecasting. *Energy Economics* 68, 77–88.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and Édouard Duchesnay (2011). Scikit-learn: Machine learning in python. *Journal of Machine Learning Research* 12(85), 2825–2830.
- Schwager, J. D. (2017). *A Complete Guide to the Futures Market: Technical Analysis, Trading Systems, Fundamental Analysis, Options, Spreads, and Trading Principles*. John Wiley & Sons.
- Smith, T. G. et al. (2017). pmdarima: Arima estimators for Python. [Online; accessed `today`].
- Uddin, G. S., R. Gençay, S. Bekiros, and M. Sahamkhadam (2019). Enhancing the predictability of crude oil markets with hybrid wavelet approaches. *Economics Letters* 182, 50–54.
- Wang, L., Z. Zhang, and J. Chen (2017). Short-term electricity price forecasting with stacked denoising autoencoders. *IEEE Transactions on Power Systems* 32(4), 2673–2681.
- Wang, Y., F. Ma, Y. Wei, and C. Wu (2016). Forecasting realized volatility in a changing world: A dynamic model averaging approach. *Journal of Banking & Finance* 64, 136–149.
- Wang, Y., Y. Wei, C. Wu, and L. Yin (2018). Oil and the short-term predictability of stock return volatility. *Journal of Empirical Finance* 47, 90–104.
- Wei, Y., J. Liu, X. Lai, and Y. Hu (2017). Which determinant is the most informative in forecasting crude oil market volatility: Fundamental, speculation, or uncertainty? *Energy Economics* 68, 141–150.
- Yin, L. and Q. Yang (2016). Predicting the oil prices: Do technical indicators help? *Energy Economics* 56, 338–350.
- Zhang, Y., F. Ma, B. Shi, and D. Huang (2018). Forecasting the prices of crude oil: An iterated combination approach. *Energy Economics* 70, 472–483.
- Zhang, Y., Y. Wei, Y. Zhang, and D. Jin (2019). Forecasting oil price volatility: Forecast combination versus shrinkage method. *Energy Economics* 80, 423–433.

## Appendix A. Data Description

### Rolling Return Strategy

We forecast the monthly volatility for futures contracts, nearest to maturity. Because each future contract must be exercised at the maturity date, to estimate continuous volatility, daily future prices are calculated based on the roll-over rule.<sup>12</sup> On each monthly roll-over day, the oil price is adjusted by the price difference between the future contract with the first nearest maturity and the future contract with the second nearest maturity. The monthly volatility of oil price is calculated by the standard deviation of daily WTI prices at each last business day of the month. A large body of literature investigated the forecasting oil price volatility using historical realized volatility. For example, Ma et al. (2017), Ma et al. (2018), and Wei et al. (2017) studied forecasting models of realized oil volatility.

### Data Preprocessing

Data needs preprocessing procedure for distinct frequency, stationarity, and integration before analysis. If the date of the publication of infrequent data is not specified, we assume that it is announced at the last moment of the month (or the quarter).<sup>13</sup> For example, the Fed Funds rate or MSCI data are monthly, but their release dates are specified. Therefore, we use announced data from the exact date that they are announced and fill missing values by using them until the next announcement date. When the announced date is not specified (such as categorical EPU data), we use them from one day after its announced month (or quarter). Each categorical data, such as oil production, consumption, and export, are country-level data. They are summed to make a global value of each categorical data. For instance, country-specific oil production data are aggregated to one global oil production.

The stationarity of all data from April 2002 to April 2024 is checked using Python's `pmdarima` package (Smith et al., 2017) except the lagged volatility variable, the mean of lagged volatility, and the dependent variable itself. If we have to take the

---

<sup>12</sup>Each contract expires on the third business day before the 25th calendar day. When the 25th calendar day is not a business day, the contract expires on the third business day before the latest business day prior to the 25th calendar day.

<sup>13</sup>Capacity utilization rate is the only weekly disclosed data, and each date of the announcement is recorded.

difference of each sample more than twice to meet the criteria of stationarity (such as adjusted Dicky-Fuller test), we set them to be 2. Even if high order differentiation may ensure the stationarity, we believe those process will omit valuable information.

## Appendix B. Training Models

Our forecasting models consist of two parts: feature selection and training. We first select features for each sample using the f-regression method in Python’s scikit-learn package (Pedregosa et al., 2011). As different sample has different characteristics, there is no guarantee that all samples use the same features. Feature selection helps models train in a reasonable time and helps prevent the notorious issue of over-fitting. After the aforementioned transformation, we fit our data to various models. Because there could be randomness during training models, scikit-learn models’ random-states are fixed to a specific number. Without this, the model will predict different values each time we train.

As the optimal predictive model may vary depending on the sample, the optimal hyperparameters may also vary. Therefore, for each training sample, different hyperparameters may be used for each model. Scikit-learn’s GridSearchCV module is an effective tool for identifying these optimal hyperparameters, and we have adopted it for this purpose.

## Appendix C. Methods of Empirical Analysis

### Diebold-Marino Test

The Diebold-Marino (DM) test is a method of comparing the performance of two or more forecasting models, introduced by Diebold and Mariano (2002) and modified by Harvey et al. (1997). The DM test compares the forecasting errors of two models and statistically verifies whether one model **overpredicts or underpredicts** compared to other models. When  $g(e_{i,t})$  is the squared-error loss or the absolute error loss of the

forecasting error  $e_{i,t} = \hat{y}_{i,t} - y_{i,t}$  for model  $i = 1, 2$ , the loss differential is defined by  $d_t = g(e_{1,t}) - g(e_{2,t})$ . The null hypothesis indicates that prediction accuracy for two models are same. The DM statistics is defined as  $DM = \frac{\bar{d}}{\sqrt{\hat{\gamma}_0/h}} \sim N(0, 1)$ , where  $\bar{d} = \sum_{t=1}^T d_t/h$  is the sample mean with  $h$ -step forward forecast,  $\bar{\gamma}_0$  is the consistent estimate of the variance for  $h\bar{d}$ .

$$H_0 : E[d_t] = 0, \quad H_1 : E[d_t] \neq 0, \quad \forall t$$