

Forecasting Cryptocurrencies Log>Returns: a LASSO-VAR and Sentiment Approach

Milos Ciganovic* Federico D’Amario†

January 9, 2024

Abstract

Cryptocurrencies have become a trendy topic recently, primarily due to their disruptive potential and reports of unprecedented returns. Furthermore, social media has garnered attention for its predictive capabilities in various fields, including financial markets and the economy. In this study, we exploit the predictive power of sentiment from Twitter and Reddit, alongside Google Trends indexes, to forecast log returns for ten cryptocurrencies, namely *Bitcoin*, *Ethereum*, *Tether*, *Binance Coin*, *Litecoin*, *Enjin Coin*, *Horizen*, *Namecoin*, *Peercoin*, and *Feathercoin*. We evaluate the performance of LASSO Vector Autoregression using daily data from January 2018 to January 2022. In a 30-day recursive forecast, we achieve a mean directional accuracy (MDA) rate of over 50%. Moreover, we observe a significant increase in forecast accuracy in terms of MDA when using sentiment and attention variables as predictors, but only for less capitalized cryptocurrencies. This improvement is not reflected in the RMSE. We also conduct a Granger causality test using post-double LASSO selection for high-dimensional VAR models. Our results suggest that social media sentiment does not Granger-cause cryptocurrencies returns.

Keywords: Cryptocurrencies, Time series analysis, Sentiment analysis, Natural Language Processing

JEL Codes: C32, C53, C55, G17

*Department of Economics and Law - Sapienza University of Rome. milos.ciganovic@uniroma1.it

†Department of Economics and Law - Sapienza University of Rome. federico.damario@uniroma1.it

1 Introduction

The rise of cryptocurrencies has been a significant phenomenon since the introduction of Bitcoin by Nakamoto (2008). Over the last decade, cryptocurrencies have experienced tremendous growth in their market capitalization and the number of kinds of coins. Several factors have contributed to this boom, including the reported unprecedented returns of cryptocurrencies in social media and journals, which led to a surge of interest among investors, both professional and non-professional. This was due, in part, to the minimal global regulation of cryptocurrencies and their perceived use for illegal trades. The potential for significant returns has also stimulated enthusiasm among investors, much like the Gold Rush in the western U.S.

Investing in cryptocurrencies is relatively easy, requiring only the downloading of an app on a smartphone. Consequently, young people represent a significant proportion of cryptocurrency investors¹. Given the young age of the cryptocurrency market, traditional news outlets are often unable to keep pace with events, leading social media to become the primary source of information for cryptocurrency investors. Micro-blogging websites like Twitter² and Reddit³ are commonly used sources for cryptocurrency information.

The high volatility and significant fluctuations in the prices of cryptocurrencies have created substantial risks for investors, resulting in heated discussions about their role in the modern economy (e.g., Corbet et al. (2019), Catalini and Gans (2020), Halaburda et al. (2020), and Auer et al. (2021)). Therefore, developing appropriate methods and models for predicting prices for digital currencies is relevant for the scientific community and financial analysts, investors, and traders. Researchers have made significant contributions to this area, including Catania et al. (2019) and Catania and Grassi (2021), who developed a dynamic model suitable for the complex dynamics of cryptocurrency time series and compared several alternative univariate and multivariate models for point and density forecasts. In addition, many studies have shown that Machine Learning algorithms, such as those presented in Hitam and Ismail (2018), Sun et al. (2020), and Miller and Kim (2021), are extremely convenient in terms of computational time and accuracy when forecasting cryptocurrency time series. Recent literature has emphasized the importance of specific factors that influence the demand for cryptocurrencies, which can aid in forecasting their prices, returns, and volatil-

¹See <https://www.investopedia.com/younger-generations-bullish-on-cryptocurrencies-5223563>.

²See <https://twitter.com/>.

³See <https://www.reddit.com/>.

ity. According to the efficient market hypothesis (EMH) Fama (1970), market prices reflect all available information, which implies that predicting stock returns should not be possible. On the other hand, considerable empirical evidences (see, e.g. Daniel et al. (2002) for a comprehensive review) show that investors' psychology drives the stock market. Specifically for Bitcoin, Ciaian et al. (2016) show that attention-driven behaviour from both investors and users affects the Bitcoin price formation in the short run. This contribution can be either positive or negative, depending on the type and the point in time of news. These findings led many researchers to adopt sentiment indexes and attention variables to improve forecast accuracy. Glenski et al. (2019) exploit the predictive power of social signals from multiple platforms (GitHub and Reddit) to forecast prices for three cryptocurrencies. They show that social signals reduce error when forecasting daily coin prices and that the language used in comments within the official communities on Reddit are the best predictors overall. Kraaijeveld and De Smedt (2020) show that Twitter sentiment has predictive power for the returns of several cryptocurrencies. Aslanidis et al. (2022) and Nasir et al. (2019) highlight the link of Google Trends with cryptocurrencies regarding their returns and volatility.

In this study, we analyze the impact of sentiment variables on cryptocurrency returns by using a novel dataset that combines a number of social media, search engine data, and volume. We apply a state-of-the-art sentiment classification technique to investigate whether sentiment measures contain predictive power for returns. To the best of our knowledge, similar sets of predictors have not been employed jointly previously. We account for the high dimensionality of the predictor variables by using a regularization technique known as the LASSO. This allows us to investigate (i) whether the variables constructed from our novel dataset can help to improve log-return forecasts using a VAR approach compared to the benchmark models; (ii) which data source and which type of sentiment or attention measure is most relevant in terms of Granger-causality in High-Dimensional VARs.

Our findings indicate that the performance of LASSO-VAR models, measured in terms of Mean Directional Accuracy (MDA), is comparable to that of benchmark models. Additionally, incorporating sentiment and attention variables as predictors significantly enhances the forecast accuracy for less capitalized cryptocurrencies. Our results do not reveal any Granger causality from sentiment indexes to cryptocurrencies returns. Instead, we observe Granger causality among all cryptocurrencies except Bitcoin, Tether, and Feathercoin, as well as from returns to the bitcoin sentiment extracted from Twitter. To further test the robustness of our approach, we conduct an additional analysis forecasting the returns of five

Nasdaq stocks, with and without the use of sentiment variables. As confirmation of our results, we do not find evidence of improved forecast accuracy when incorporating sentiment variables for these stocks.

The remainder of the paper is organized as follows. Section 2 deals with data collection, describing the dataset, its sources, and strategies to construct sentiment indexes. Section 3 describes the modelling strategy, the estimation, forecasting methods, and the metrics used for the comparative evaluation of the out-of-sample model predictions. Section 4 summarizes some selected results. Section 5 presents robustness checks. Section 6 concludes.

2 Data Collection

The present study utilizes multiple data sources to predict cryptocurrency returns and study the influence and causality of sentiment variables on target variables. Specifically, we obtain daily data of ten cryptocurrencies from January 2018 to January 2022, as well as Google Trends and sentiment indexes from Twitter and Reddit for the same period. To further enhance our dataset, we include volume data for each of the considered cryptocurrencies.

2.1 Cryptocurrency Data

The cryptocurrencies examined in this study are listed in Table 1, which presents them in order of Market Capitalization (MC) as of January 2022. The data utilized in this investigation were obtained from [finance.yahoo.com](https://it.finance.yahoo.com/cryptocurrencies/)⁴.

Many reasons brought us to choose this set of cryptocurrencies. First of all, cryptocurrencies emerge and disappear continually, while our selected ten currencies have been publicly-traded consecutively. Moreover, all the currencies chosen have been created with a defined purpose representing innovative projects which brought development, progress, or value to the blockchain technology that Bitcoin had implemented. On the other hand, our sample includes three tier currencies as in Gandal and Halaburda (2016). *Bitcoin*, *Ethereum*, *Tether*, *BinanceCoin*, whose market capitalizations stay in the world’s top five, are “top-tier” cryptocurrencies. *Litecoin*, *EnjinCoin*, *Horizen*, representing “middle cryptocurrencies” in market capitalization. *Namecoin*, *Peercoin*, and *Feathercoin* are representative “minor cryptocurrencies” according to market capitalization. We include Tether (USDT) in our sample for a specific reason. We know, indeed, that it is a blockchain-based cryptocurrency

⁴See: <https://it.finance.yahoo.com/cryptocurrencies/>

Table 1: 10 Cryptocurrencies and their symbols, market capitalization (MCs) (as of 31 January 2022).

Cryptocurrency	Symbol	MC	Rank by MC
Bitcoin	BTC	\$702,864,225,136	1
Ethereum	ETH	\$295,905,148,931	2
Tether	USDT	\$78,188,468,450	3
Binance Coin	BNB	\$63,930,448,963	4
Litecoin	LTC	\$7,479,561,631	21
Enjin Coin	ENJ	\$1,434,490,287	60
Horizen	ZEN	\$505,212,977	120
Namecoin	NMC	\$24,472,607	733
Peercoin	PPC	\$13,809,105	837
Feathercoin	FTC	\$1,927,251	1515

whose tokens in circulation are backed by an equivalent amount of U.S. dollars, making it a stablecoin with a price pegged to USD \$1.00, which leads this currency to be very low volatile. Tether was designed to build the necessary bridge between fiat currencies and cryptocurrencies and offer users stability, transparency, and minimal transaction charges. We decided to include it as a counterfactual to understand whether sentiment indexes can help the prediction of stablecoin currencies. We compute the log returns and include them in our sample. Table 2 provides several summary statistics.

Table 2: Log-returns summary statistics for the ten cryptocurrencies during the period 1 January 2018 to 31 January 2022

	BTC-USD	ETH-USD	USDT-USD	BNB-USD	LTC-USD	ENJ-USD	ZEN-USD	NMC-USD	PPC-USD	FTC-USD
Mean	0.001	0.001	0	0.003	0	0.002	0	-0.001	-0.001	-0.003
Median	0.001	0.001	0	0.001	0	-0.001	-0.001	0.001	-0.001	-0.005
Min	-0.465	-0.551	-0.053	-0.543	-0.449	-0.624	-0.546	-1.16	-0.665	-0.474
Max	0.172	0.231	0.053	0.529	0.291	0.768	0.38	0.75	0.567	0.409
Range	0.637	0.781	0.106	1.072	0.74	1.392	0.926	1.91	1.232	0.883
Skew	-1.147	-1.099	0.3	0.3	-0.608	1.132	-0.233	-0.869	-0.174	-0.174
kurtosis	14.042	10.795	34.34	15.579	7.947	15.508	6.268	19.999	13.593	4.575
ADF	-11.0753	-11.1057	-14.186	-11.1936	-11.3664	-11.201	-11.0829	-12.9631	-12.9892	-11.5868

Notes: all ADF statistics are stationary at 1% level

2.2 Google Trends Data

We collected forty three google trends searches (Table 9 in the Appendix reports the list of all the Google trends collected).

Figure 1: Google trends words collected



Google Trends is a search trend feature that shows how frequently a given search term is entered into Google’s search engine relative to the site’s total search volume over a given period. Google Trends are available at a daily frequency only if the selected period is less than nine months. If the time frame is between nine months and five years, weekly data are provided; if it is longer than five years, data are monthly. A trivial solution like querying the data month by month and then tying it together will not work in this case because Google Trends assess interest in relative values within the given time. It means that for a given keyword and month, Google Trend will estimate the interest identically, with a local minimum of 0 and a local maximum of 100, for events in one month even if they had twice as many searches than in the other. To get proper daily estimates, we proceed as follows:

- 1 - Query daily estimates for each month in the specified time frame;
- 2 - Queries monthly data for the whole time frame;
- 3 - Multiply daily estimates for each month from step 1 by its weight from step 2.

2.3 Twitter and Reddit Sentiment Indexes

With regard to sentiment indexes, we use the Twitter API v2 ⁵ to download tweets using the name of each cryptocurrency from our sample as a search parameter. We end up with 13,195,084 tweets from 2,253,469 unique users from January 2018 to January 2022. Each tweet’s query provides the user with the timestamp (UTC+0) and the text containing a maximum of 280 characters. Concerning Reddit data, we use the Pushshift Reddit API

⁵See <https://developer.twitter.com/en/docs/twitter-api>

Baumgartner et al. (2020) to collect comments under the official subreddit of each cryptocurrency included in our sample. We gather a total of 4,406,897 comments from 420,024 unique users in the same time frame specified above. We are aware that increasing the daily quantity of tweets and comments collected would enhance the quality and results of our study. However, we believe that our sample is enough to get the representative social sentiment we are looking for.

The sentiment time series are created using Valence Aware Dictionary and Sentiment Reasoner (VADER) algorithm Hutto and Gilbert (2014). It uses a list of lexical features, optimized for social media texts, (words, punctuation, and emoticons) labelled as positive or negative according to their semantic orientation to calculate the text sentiment. Vader sentiment returns the probability of a given input sentence to be positive, negative, and neutral. In particular, every tweet and comment are passed to the algorithm. Therefore, we retrieve a Valence score for each of them which is measured on a scale from -4 to +4, where -4 stands for the most “Negative” sentiment and +4 for the most “Positive” sentiment. Midpoint 0 represents a “Neutral” Sentiment. The unidimensional measure the algorithm can retrieve is called the “Compound Score”. It is computed by summing the valence scores of each word in the lexicon, adjusted according to the rules, and then normalized to be between -1 (most extreme negative) and +1 (most extreme positive):

$$C = \frac{x}{\sqrt{x^2 + \alpha}} \quad (1)$$

Where x is the sum of the valence score of constituent words, and α is a normalization constant (the default value is 15). Finally, we obtain the compound score for each input processed and aggregate them into daily frequency. In particular, for Twitter, we get the daily measure of sentiment (DS) by weighting each Tweet score by its number of retweets.

$$DS = \frac{\sum_{i=1}^I r_i C_i}{\sum_{i=1}^I r_i} \quad (2)$$

Where I is the total number of Tweets in a day, r_i is the number of retweets. It is essential to say that we increased the number of retweets by one. This procedure allows us to perform the weighted average in (2) for those days where Tweets had zero retweets. At the same time, weighting Tweets with zero retweets is like assuming they are shared only by the author and nobody else. With this procedure, we seek to give more importance to those Tweets that most people share among their community. In Table 10 we report the cross-correlation

coefficients between cryptocurrencies and their specific sentiment from Twitter and Reddit. Most average cryptocurrency values at time $t = 0$ are significantly correlated with lagged observations of their specific sentiment up to one week before. This is crucial evidence of how influential sentiment variables are in our models to increase predictive accuracy.

2.4 Volume

We complete our dataset with cryptocurrencies volume. The dynamic volume-return could bring valuable information which can be used for price predictions or trading strategies. Generally speaking, the literature concerning the relationship between volume and returns suggests that there is a positive correlation between them (see, e.g. Jain and Joh (1988) and Llorente et al. (2002)). Concerning cryptocurrencies, there is a growing literature studying the volume-return relation. Balcilar et al. (2017) employ a non-parametric causality-in-quantiles test to assess the causality relation between trading volumes and bitcoin returns and volatility. They show that volume can predict returns, except during bull and bear periods in the Bitcoin market. Zhang et al. (2018) investigate the return-volume relationship of the Bitcoin market based on multifractal detrended cross-correlation. They found that Bitcoin exhibits a non-linear dependent relationship in return-volume with cross-correlations of return-volume showing anti-persistent behaviour. Another, Naeem et al. (2020) explore extreme return-volumes dependence among the highest capitalized cryptocurrencies using the Copula approach. They discovered that coefficients of lower tail dependence are significant for Bitcoin, Ripple, and Litecoin, which means that low volumes follow low returns. Lower tail dependence for the return-volume relationship is stronger than the upper tail dependence for Bitcoin, Ripple, and Litecoin. Moreover, for negative return-volume, left tail dependence coefficients are significant for Ripple and Litecoin, which means that low volumes follow high returns for Ripple and Litecoin. Finally, in a recent work, Chan et al. (2022) investigate the extreme dependence and correlation between high-frequency cryptocurrency returns and transaction volumes, at the extreme tails associated with booms and busts in the cryptocurrency markets. Applying an extreme value theory approach, they found a weak positive correlation between return and volume at the tails.

3 Methods

3.1 LASSO-VAR

The LASSO is a method for automatic variable selection and parameters shrinkage. It can be used to select the most informative predictors of a target variable Y from a set of variables and parameters, possibly larger than sample information, virtually making high-dimensional modelling and forecasting feasible for any degree of model dimension and complexity. The LASSO has been initially developed for a single equation setting by Tibshirani (1996). The LASSO approaches curve fitting as a quadratic programming problem, where the objective function penalizes the total size of the regression coefficients based on the value of a tuning parameter, λ . In doing so, the LASSO can drive the coefficients of irrelevant variables to zero, thus performing the automatic variable selection. The strength of the penalty must be tuned. The stronger the penalty, the higher the number of coefficients that are shrunk to zero. The model is thus forced to select only the most important predictors, i.e. those with the highest contribution to the prediction of the target variable.

Let $\{y_t\}_{t=1}^T$ be a K dimensional multiple time series process generated by VAR process of order p , denoted as VAR(p):

$$\begin{aligned} y_t &= A_1 y_{t-1} + \dots + A_p y_{t-p} + u_t, \\ u_t &\sim \mathcal{N}(0, \Sigma_u). \end{aligned} \tag{3}$$

In our analysis, we fix the maximum order of lag p to fourteen days, as suggested by the Bayesian Information Criterion, calculated over the full sample, i.e. including data up to the end of January 2022. Notice that, since we standardize data before modelling, the K -dimensional intercept vector is not considered in the VAR. Each A_i is a $K \times K$ matrix of coefficients for the endogenous variables, and $u_t \stackrel{\text{wn}}{\sim} (0, \Sigma_u)$ is the vector of reduced-form errors.

For notation convenience we will introduce a compact matrix representation of (3)

$$\begin{aligned}
\mathbf{Y} &= [y_{p+1}, \dots, y_T] & (K \times T); \\
\mathbf{Z}_t &= [y'_p, \dots, y'_{t-p}] & [1 \times Kp]; \\
\mathbf{Z} &= [\mathbf{Z}_p; \dots; \mathbf{Z}_{T-p}] & (T \times Kp)'; \\
\mathbf{A} &= [A_1, \dots, A_p] & (K \times Kp); \\
\mathbf{U} &= [u_1, \dots, u_T] & (K \times T)
\end{aligned}$$

We can express the VAR in (3) as

$$\mathbf{Y} = \mathbf{AZ} + \mathbf{U}, \quad (4)$$

with $\mathbf{U} \sim \mathcal{N}(0, \mathbf{I}_t \otimes \boldsymbol{\Sigma}_u)$.

The LASSO objective function is minimized as follows:

$$\hat{\mathbf{A}}(\lambda) = \arg \min_{\mathbf{A}} \frac{1}{T} \|\mathbf{AZ} - \mathbf{Y}\|_2^2 + \lambda \|\mathbf{A}\|_1, \quad (5)$$

where λ is the shrinkage parameter.

Together with the homoskedastic framework described above, we decide to account for heteroskedasticity estimating the model through Feasible GLS estimator assuming an $AR(1)$ process for the errors. In particular, we take the estimated residuals of the homoskedastic case and we regress them to their lagged value obtaining the estimate for the autocorrelation parameter.

For each k^{th} equation in the VAR in (3) we will have:

$$\hat{u}_{k,t} = \rho_k \hat{u}_{k,t-1} + \epsilon_{k,t} \quad (6)$$

where $|\rho_k| < 1$ and $\epsilon_{k,t}$ is white noise with variance σ_k^2 . Then

$$E(u_k u'_k) = \frac{\sigma_k^2}{1 - \rho_k^2} \begin{bmatrix} 1 & \rho_k & \rho_k^2 & \dots & \rho_k^{T-1} \\ \rho_k & 1 & \rho_k & \dots & \rho_k^{T-2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \rho_k^{T-1} & \rho_k^{T-2} & \rho_k^{T-3} & \dots & 1 \end{bmatrix} = \sigma_k^2 \mathbf{T}_k \quad (7)$$

Since \mathbf{T}_k is symmetric and positive definite, there exists a nonsingular $(T \times T)$ matrix \mathbf{L}_k such that:

$$\mathbf{T}_k^{-1} = \mathbf{L}'_k \mathbf{L}_k \quad (8)$$

We can the rewrite the VAR in (4) as:

$$\mathbf{Y}\mathbf{L}_k = \mathbf{A}\mathbf{Z}\mathbf{L}_k + \mathbf{U}\mathbf{L}_k \equiv \tilde{\mathbf{Y}} = \mathbf{A}\tilde{\mathbf{Z}} + \tilde{\mathbf{U}} \quad (9)$$

with $\tilde{\mathbf{U}} \sim \mathcal{N}(0, \mathbf{I}_t \otimes \boldsymbol{\Omega}_u)$.

The Feasible GLS LASSO objective function is minimized as follows:

$$\hat{\mathbf{A}}(\lambda) = \arg \min_{\mathbf{A}} \frac{1}{T} \|\mathbf{A}\tilde{\mathbf{Z}} - \tilde{\mathbf{Y}}\|_2^2 + \lambda \|\mathbf{A}\|_1, \quad (10)$$

By applying the LASSO procedure, we seek to obtain a sparse structure for the coefficient matrices. The optimization problem is solved by applying a coordinate descent numerical procedure, as explained in Kim et al. (2007) and Friedman et al. (2010).

3.2 Calibrating the LASSO-VAR through time series cross-validation

As immediately evident from (10), lambda (λ) is the most important parameter in the LASSO framework. The selection of the best predicting model depends on its calibration, which should not be sample-specific. A cross-validation stage is thus employed to get the “optimal” value for λ .

In this respect, we emply an expanding window (more precisely, an “anchored walk forward”) approach to cross-validation, which is the standard practice in many analyses considering time series modelling. In practice, the data is divided into a training and a test sample. The test set is being held for final evaluation, whereas the training set is further split into three subsets.

An example of 3-split time series cross-validation method is depicted in Figure 2.

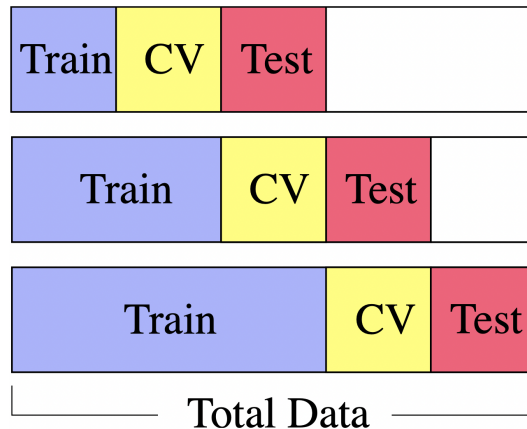


Figure 2: Graphical example of the 3-split time series cross validation

Time series cross-validation essentially considers time dependence, such that the training set is always built considering observations prior to the test set. The anchored walk-forward cross-validation method implies a gradually expanding training set, pushing forward a fixed dimension test set. In practice, we chose the shrinkage intensity parameter (λ) at each recursive forecast step, which optimizes the results of out-of-sample forecast in terms Root Mean Squared Error. We decide for a 5-split time series cross validation considering forecast performances and computational time.

3.3 Evaluation of the forecasting performances

After estimating the sparse model resulting from the cross-validated LASSO-VAR, we proceed to forecast. We use the sample from 01/01/2018–12/31/2021 to obtain initial parameter estimates for all models, which are then used to predict outcomes from 01/01/2022 ($h = 1$) to 01/04/2022 ($h = 4$). In the next period, we include data for 01/01/2022 in the estimation sample and use the resulting estimates to predict the outcomes from 01/02/2022 to 01/06/2022. We proceed recursively in this fashion until 01/31/2022, thus generating a time series of point forecasts. In (11) we show the standard procedure to obtain out-of-sample forecasts when a VAR is assumed to be the data generating process.

$$y_{T+h|T} = c + A_1 y_{T+h-1|T} + \dots + A_p y_{T+h-p|T}. \quad (11)$$

In order to evaluate the forecasting model performances, we adopt the Root Mean Squared Error (RMSE) and the Mean Directional Accuracy (MDA):

$$RMSE = \sqrt{\frac{1}{H} \sum_{\tau=\underline{t}}^{\bar{t}-h} e_{i,\tau+h}^2}, \quad (12)$$

$$MDA = \frac{1}{H} \sum_{\tau=\underline{t}}^{\bar{t}-h} \text{sign}(y_{i,\tau+h} - y_{i,\tau+h-1}) == \text{sign}(\hat{y}_{i,\tau+h} - \hat{y}_{i,\tau+h-1}). \quad (13)$$

We evaluate the performances of our selected models by comparing the results with the two multivariate high-dimensional time series forecasting model alternatives in discrete time. Specifically, we consider the Large Bayesian VAR (LBVAR) introduced by Bańbura et al. (2010) and Factor Augmented VAR (FAVAR) proposed by Bernanke et al. (2005).

The latter model, in particular, is estimated using specific factors for Google Trends, Twitter

and Reddit Sentiments and, Volumes. We created the factor series summing the orthogonal principal components so that we could achieve at least the 50% of explained variance for each group. Finally, in order to test the statistical significance of differences in point forecasts, we consider pairwise tests of equal predictive accuracy (henceforth, EPA; Diebold and Mariano (2002); West (1996)) in terms of RMSE. All EPA tests we conduct are based on a two sided test with the null hypothesis being the LBVAR and FAVAR benchmarks. We use standard normal critical values. Based on simulation evidence in Clark and McCracken (2013), when computing the variance estimator which enters the test statistic we rely on serial correlation robust standard errors and incorporate the finite sample correction due to Harvey et al. (1997). In Table 6, we use ***, ** and * to denote results which are significant at the 1%, 5% and 10% levels, respectively, in favor of the model listed at the top of each column compared with LBVAR. We use ***, ** and * when we compare our models with FAVAR.

3.4 High-dimensional Granger causality test

The concept of Granger causality captures the predictability given a particular information set Granger (1969, 1980). If the addition of variable I to the given information set Ξ alters the conditional distribution of another variable J , and both I and Ξ are observed prior to J , then I improves the predictability of J and is said to Granger cause J with respect to Ξ . Granger (1969) considers Ξ in a theoretical and non-practicable way as all the information available in the universe. The choice of the information set plays thus a crucial role in determining (non-)Granger causality. Cases of spurious Granger causality from I to J may arise when both I and J are Granger caused by Q , but Q is omitted from Ξ . High-dimensional models may avoid possible spurious causality due to omitted variables. Of course, conditioning on so many variables leads to the curse of dimensionality, rendering many standard statistical techniques invalid. To overcome the problem mentioned above, regularized estimation procedures as the LASSO can be used. As described in Hecq et al. (2019), one might be tempted to perform the LASSO as in (10) on (14)⁶:

$$y_J = Z^\otimes \beta + u_J = Z_{GC}^\otimes \beta_{GC} + Z_{-GC}^\otimes \beta_{-GC} + u_J, \quad (14)$$

⁶Where $y_J = \text{vec}(\mathbf{Y}_J)$ denote the $N_J \times 1$ stacked vector containing all the observations of the variables in J . Similarly $u_J = \text{vec}(\mathbf{U}_J)$, $Z^\otimes = \mathbf{I}_{N_J} \otimes \mathbf{Z}$ and $\beta = \text{vec}(\mathbf{A})$. GC and $-GC$ stand respectively for chosen possible Granger causing variables and remaining variables.

testing the parameter significance of the selected “Granger-causing” variables. However, this procedure does not consider that the final and selected model is random and function of the data. The randomness in the selection step means that post-selection estimators do not converge uniformly to a normal distribution, as the potential omitted variable bias from omitting weak, but still relevant, variables in the selection step is too large to perform uniformly valid inference. Many authors tried to cope with this issue (see Hecq et al. (2019) for a comprehensive review). We adopt the procedure proposed by Hecq et al. (2019) who specifically implemented a post-double-selection procedure, initially developed by Belloni et al. (2014), in a VAR context. The idea behind their approach is that Ξ is made of variables of interest, possible Granger-causing variables and the remaining variables. Their idea is to regress both the variable of interest and possible Granger-causing variables on the remaining variables. These regressions could be high-dimensional and cannot be estimated by least square. The remaining variables retained are those whose parameter is significant in at least one of the regressions mentioned above. In this case, the omitted variable bias will only occur if the LASSO fails to select a relevant variable in both regressions simultaneously. As the probability for this to occur decreases quadratically, this is negligible asymptotically and allows for valid inference.

In our application we test the (non-)Granger Causality for all the variables in our dataset considering all the possible combinations. Furthermore, since the time series length of our variables is much higher in magnitude compared to the number of covariates, we do not need to set any bound on the penalty to ensure that in each regression a maximum of variables are selected. In this case, we tune λ selecting the one which minimizes the Bayesian information criterion (BIC).

4 Results and analysis

In this study, we undertake a comparison of our proposed methodologies with the alternative models that were previously introduced in Section 3.3 evaluating the performances removing sentiment and attention variables. To this end, we present the MDA statistics for the return of each cryptocurrency in Table 3, which summarizes the forecasting accuracy of each model. The table is divided into four sections, each focusing on a different forecast horizon, with the rows pertaining to individual cryptocurrencies. Our analysis reveals that FGLS LASSO-VAR and LBVAR exhibit superior performance in terms of MDA. Notably, the heteroscedastic VAR model achieves better results for longer forecast horizons. Additionally, our findings indicate significant variations in the forecasting ability of cryptocurrencies that belong to the three tiers classified in Section 2. Specifically, we observe that minor cryptocurrencies exhibit comparatively lower MDA results when compared to other digital currencies.

Table 3: The table shows the Mean Directional Accuracy of model i for variable j computed as $MDA = \frac{1}{N} \sum_{t=\underline{t}}^{\bar{t}} \text{sign}(y_{t+1,i,j} - y_{t,i,j}) == \text{sign}(f_{t+1,i,j} - f_{t,i,j})$. Where $N = 29$ are the directions for 30 days forecast. y is the actual value and f is the forecast. Bold values are the best results.

Variable	FGLS LASSO-VAR	OLS LASSO-VAR	LBVAR	FAVAR	Variable	FGLS LASSO-VAR	OLS LASSO-VAR	LBVAR	FAVAR
h=1					h=2				
BTC-USD	52.00%	48.00%	60.00%	56.00%	BTC-USD	50.00%	57.69%	65.38%	50.00%
ETH-USD	40.00%	52.00%	52.00%	44.00%	ETH-USD	57.69%	53.85%	61.54%	50.00%
USDT-USD	60.00%	64.00%	60.00%	60.00%	USDT-USD	73.08%	65.38%	46.15%	69.23%
BNB-USD	56.00%	48.00%	52.00%	36.00%	BNB-USD	65.38%	61.54%	61.54%	30.77%
LTC-USD	48.00%	52.00%	60.00%	52.00%	LTC-USD	53.85%	57.69%	65.38%	50.00%
ENJ-USD	52.00%	56.00%	40.00%	44.00%	ENJ-USD	53.85%	69.23%	50.00%	42.31%
ZEN-USD	64.00%	56.00%	56.00%	52.00%	ZEN-USD	65.38%	61.54%	61.54%	57.69%
NMC-USD	48.00%	52.00%	72.00%	64.00%	NMC-USD	53.85%	57.69%	65.38%	61.54%
PPC-USD	36.00%	24.00%	48.00%	40.00%	PPC-USD	34.62%	38.46%	50.00%	46.15%
FTC-USD	40.00%	40.00%	76.00%	80.00%	FTC-USD	50.00%	42.31%	69.23%	69.23%
h=3					h=4				
BTC-USD	51.85%	51.85%	62.96%	44.44%	BTC-USD	53.57%	50.00%	53.57%	46.43%
ETH-USD	51.85%	55.56%	59.26%	55.56%	ETH-USD	53.57%	57.14%	60.71%	57.14%
USDT-USD	81.48%	70.37%	48.15%	59.26%	USDT-USD	75.00%	67.86%	42.86%	60.71%
BNB-USD	66.67%	66.67%	51.85%	25.93%	BNB-USD	60.71%	53.57%	42.86%	35.71%
LTC-USD	44.44%	51.85%	59.26%	48.15%	LTC-USD	42.86%	46.43%	64.29%	42.86%
ENJ-USD	62.96%	55.56%	51.85%	44.44%	ENJ-USD	60.71%	53.57%	50.00%	50.00%
ZEN-USD	62.96%	59.26%	59.26%	55.56%	ZEN-USD	64.29%	53.57%	67.86%	53.57%
NMC-USD	51.85%	55.56%	59.26%	62.96%	NMC-USD	50.00%	50.00%	67.86%	71.43%
PPC-USD	44.44%	44.44%	48.15%	51.85%	PPC-USD	42.86%	46.43%	50.00%	57.14%
FTC-USD	48.15%	51.85%	66.67%	70.37%	FTC-USD	46.43%	50.00%	67.86%	75.00%

Table 4 presents the average results of each model, including the heteroscedastic model without sentiment and search engine data. Our analysis indicates that, on average, LASSO-VAR models exhibit similar directional accuracy to the LBVAR benchmark model. However, it is worth noting that the model estimated without sentiment and Google trends data

experiences a significant decrease of 16 percentage points in terms of MDA performance when compared to the full model.

Table 4: Average MDA scores obtained from each model tested.

	FGLS LASSO-VAR	OLS LASSO-VAR	FGLS LASSO-VAR (without sentiment and gtrends)	LBVAR	FAVAR
MDA h=1	51.85%	49.20%	48.00%	57.60%	52.80%
MDA h=2	55.77%	56.54%	48.46%	59.62%	52.69%
MDA h=3	56.67%	56.30%	49.63%	56.67%	51.85%
MDA h=4	55.00%	52.86%	46.79%	56.79%	55.00%

At this juncture, we seek to ascertain which cryptocurrencies' forecasts are most impacted by the exclusion of sentiment and attention variables from the model. To this end, we present the average MDA results for each cryptocurrency computed over the four forecast horizons for the FGLS LASSO-VAR model, with and without sentiments and Google Trends, in Table 5. Our analysis reveals that the forecasts for middle and minor cryptocurrencies are significantly affected by the exclusion of sentiment variables. This finding is noteworthy as it suggests that stronger and more established currencies are essentially immune to the influence of sentiment variables, whereas the price movements of minor cryptocurrencies may be swayed by sentiments, perhaps due to their nascent stage or relatively weaker ideas and projects.

Table 5: Average MDA scores computed over the forecast horizons for each cryptocurrencies.

	FGLS LASSO-VAR	FGLS LASSO-VAR (without sentiment and gtrends)
BTC-USD	51.86%	51.93%
ETH-USD	50.78%	49.22%
USDT-USD	72.39%	66.94%
BNB-USD	62.19%	55.74%
LTC-USD	47.29%	46.18%
ENJ-USD	57.38%	46.26%
ZEN-USD	64.16%	56.66%
NMC-USD	50.92%	42.55%
PPC-USD	39.48%	25.21%
FTC-USD	46.14%	41.51%

Table 6 displays the forecast performances in terms of root mean square error (RMSE). Our findings indicate that LBVAR model produces the most accurate results in absolute terms. However, these results are not statistically different from those obtained with OLS and FGLS VARs, as indicated by the equal predictive ability (EPA) tests. This outcome could be attributed to the ability of the Bayesian model to set a general and specific shrinkage intensity for each parameter, leading to more precise estimates. In contrast, both LASSO regularized models exhibit significantly better performance than FAVAR, as evidenced by the EPA tests.

Table 6: The table shows the RMSE of model i computed as $RMSE_{ij} = \sqrt{\sum_{\tau=\underline{t}}^{\bar{t}-h} e_{i,j,\tau+h}^2}$ where $e_{i,j,\tau+h}^2$ are the squared forecast error of variable j at time τ and forecast horizon $h = 1, \dots, 4$ generated by model i . \bar{t} and \underline{t} denote the start and the end of the out-of-sample period. All forecast are generated out-of-sample using recursive estimate of the model, with $\underline{t} = 01/01/2022$ and $\bar{t} = 01/31/2022$. Bold number represent the best results in a row. * significance at the 10% level; ** significance at the 5% level; *** significance at the 1% level

Variable	FGLS LASSO-VAR	OLS LASSO-VAR	LBVAR	FAVAR	Variable	FGLS LASSO-VAR	OLS LASSO-VAR	LBVAR	FAVAR
h=1					h=2				
BTC-USD	0.0423***	0.0399***	0.0324	0.1866	BTC-USD	0.0418**	0.0410**	0.0309	0.1490
ETH-USD	0.0558**	0.0584**	0.0501	0.2333	ETH-USD	0.0565**	0.0585**	0.0482	0.1862
USDT-USD	0.0043***	0.0045***	0.0007	0.1996	USDT-USD	0.0039***	0.0039***	0.0007	0.2052
BNB-USD	0.0773***	0.0784***	0.0513	0.2074	BNB-USD	0.0731***	0.0757***	0.0498	0.1824
LTC-USD	0.0610***	0.0589***	0.0458	0.2152	LTC-USD	0.0595***	0.0567***	0.0449	0.1693
ENJ-USD	0.0974***	0.0958***	0.0774	0.2902	ENJ-USD	0.0914***	0.0906***	0.0745	0.2253
ZEN-USD	0.0832***	0.0846***	0.0641	0.2872	ZEN-USD	0.0791***	0.0798***	0.0625	0.2247
NMC-USD	0.1210***	0.1209***	0.0628	0.2395	NMC-USD	0.1094**	0.1101**	0.0632	0.1830
PPC-USD	0.0763***	0.0733***	0.0447	0.1941	PPC-USD	0.0674***	0.0669***	0.0465	0.1610
FTC-USD	0.0993***	0.1000***	0.0690	0.2296	FTC-USD	0.0980**	0.0976**	0.0700	0.1817
h=3					h=4				
BTC-USD	0.0397**	0.0386**	0.0301	0.1381	BTC-USD	0.0408***	0.0392***	0.0294	0.1068
ETH-USD	0.0542*	0.0548*	0.0471	0.1642	ETH-USD	0.0537*	0.0538*	0.0462	0.1210
USDT-USD	0.0031***	0.0032***	0.0006	0.1963	USDT-USD	0.0035***	0.0037***	0.0006	0.1376
BNB-USD	0.0687***	0.0689***	0.0492	0.1719	BNB-USD	0.0673***	0.0670***	0.0479	0.1290
LTC-USD	0.0608***	0.0587***	0.0441	0.1499	LTC-USD	0.0596***	0.0571***	0.0430	0.1157
ENJ-USD	0.0882***	0.0861***	0.0723	0.2015	ENJ-USD	0.0870***	0.0844***	0.0702	0.1630
ZEN-USD	0.0753***	0.0728***	0.0606	0.2031	ZEN-USD	0.0734***	0.0712***	0.0583	0.1584
NMC-USD	0.1025*	0.1014*	0.0632	0.1611	NMC-USD	0.1008	0.0993	0.0616	0.1289
PPC-USD	0.0628***	0.063286***	0.0470	0.1349	PPC-USD	0.0624***	0.0625***	0.0499	0.1222
FTC-USD	0.0952**	0.092383**	0.0676	0.1600	FTC-USD	0.0958**	0.0925	0.0675	0.1251

Lastly, Table 7 demonstrates that sentiment variables do not enhance the forecast accuracy in terms of RMSE. This result could be attributed to the intrinsic nature of sentiment indexes. As indexes, sentiments do not furnish the model with precise information concerning the actual out-of-sample value of the returns. Instead, they provide directional information about our target variables. This translates into better performance in terms of MDA, as evidenced in Table 4, but not in terms of RMSE.

Table 7: Average RMSE scores obtained from each model tested.

	FGLS LASSO-VAR	OLS LASSO-VAR	FGLS LASSO-VAR (without sentiment and gtrends)	LBVAR	FAVAR
RMSE h=1	0.0719	0.0715	0.0774	0.0498	0.2283
RMSE h=2	0.0681	0.0681	0.0754	0.0491	0.1868
RMSE h=3	0.0651	0.0641	0.0710	0.0482	0.1681
RMSE h=4	0.0645	0.0631	0.0693	0.0475	0.1308

With regards to Granger causality, we present networks of returns, sentiment, and volume variables in Figures 3, 4, and 5, respectively. The arrows in the figures represent the direction of Granger causality when the null hypothesis is rejected with p -values of LM test < 0.01 . The p -values for each combination tested are reported in Appendix Tables: 11,12,13, and 14⁷.

Based on Figure 3, we can conclude that there is no Granger causality between Twitter sentiments and cryptocurrencies returns. However, returns seem to cause each other, and most of them cause the Bitcoin sentiment extracted from Twitter. From this finding, it is reasonable to interpret Bitcoin Twitter sentiment as a general index for all cryptocurrencies markets. It is also interesting to note that the largest and smallest cryptocurrencies in market capitalization are outside the cluster formed by all the other currencies.

⁷For simplicity, we omitted some combinations when testing Granger causality from Google trends. However, the p -values reported were computed considering search engine data as "other variables".

Cryptocurrencies – Twitter Sentimet Network

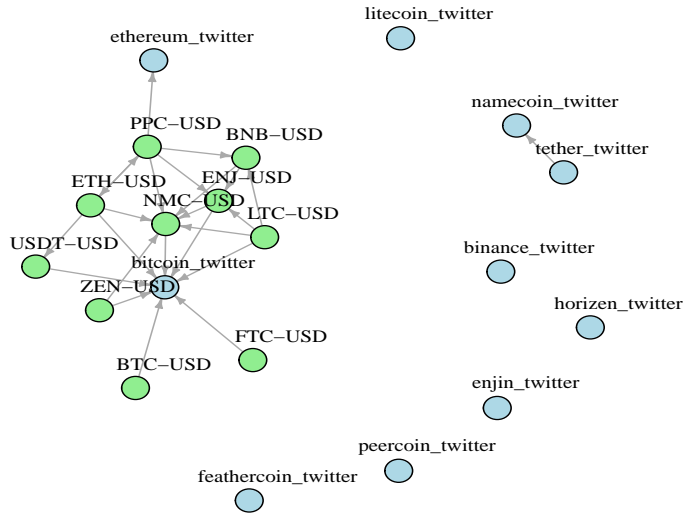


Figure 3: The figure shows a network between Twitter sentiment indexes and returns. The direction of the arrows represents the direction of the Granger causality

Same results are obtained focusing on the relation between Reddit sentiments and returns. Figure 4 shows, in fact, that sentiments obtained from Reddit are not Granger causing their respective cryptocurrency. Also, in this case, Bitcoin, Feathercoin, and the stablecoin Tether are outside the cluster created by the other currencies. It is interesting to focus on the cluster formed by several sentiments around the Bitcoin one. This finding may describe that the euphoria in specific currencies Sub-reddit is linked to the users' feelings expressed in the Bitcoin Sub-reddit, which pushes the others.

Cryptocurrencies – Reddit Sentiment Network

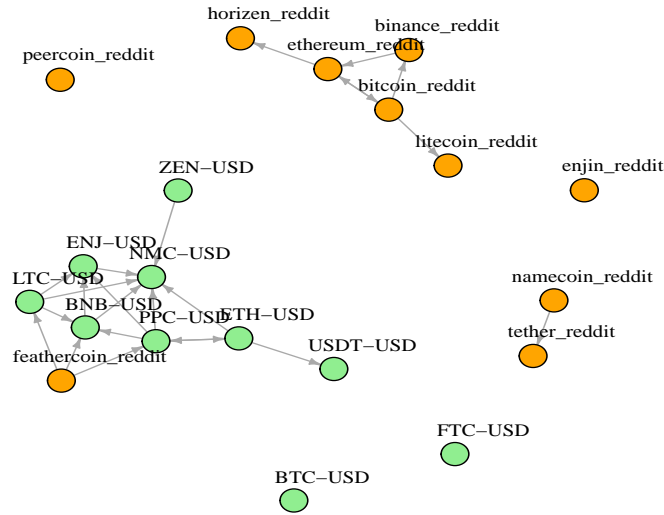


Figure 4: The figure shows a network between Reddit sentiment indexes and returns. The direction of the arrows represents the direction of the Granger causality

Figure 5 shows the Granger causality network between returns and volume.

Cryptocurrencies – Volume Network

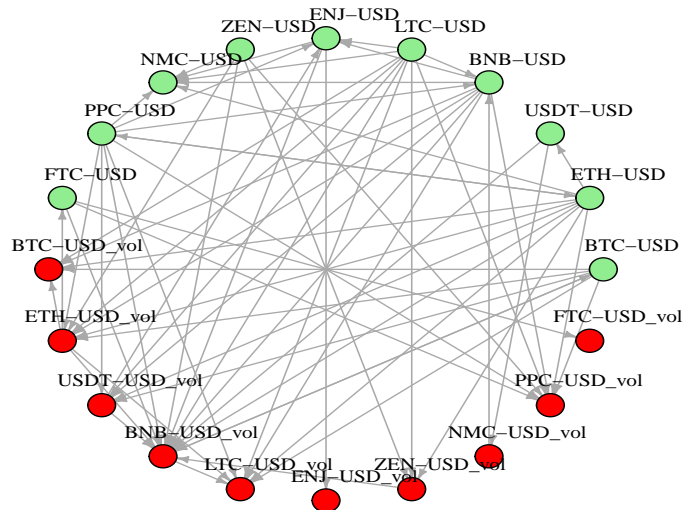


Figure 5: The figure shows a network between Volume and returns. The direction of the arrows represents the direction of the Granger causality

In this case, we can see that returns are Granger causing currencies volume. These results suggest that movements in cryptocurrency returns may influence the market of the others, primarily via demand shocks. Finally no causality is found from Google trends and

cryptocurrencies return.

5 Robustness

This section aims to investigate whether the observed positive impact of sentiment and attention variables on forecasting accuracy is unique to cryptocurrencies or if it can enhance the predictive performance of regular stocks time series. To this end, we have constructed a dataset in a similar fashion to that described in Section 2, utilizing the following stocks: *Tesla Inc. (TSLA)*, *Intel Corporation (INTC)*, *Apple Inc.(AAPL)*, *Acrivision (ATVI)*, *Starbucks Corporation (SBUX)*. Given that stocks log returns are often subject to heteroskedasticity, we employ the Feasible GLS LASSO VAR methodology to account for this characteristic. We assess the effectiveness of incorporating sentiment indexes into the forecasting model by comparing the predictive performance of models with and without these variables. The findings of this analysis are presented in Table 8.

Table 8: Average MDA and RMSE scores computed over the forecast horizons for each Stocks.

	FGLS LASSO-VAR		FGLS LASSO-VAR (without sentiment and gtrends)	
	RMSE	MDA	RMSE	MDA
TSLA	0.0642	44.23%	0.0575	51.06%
INTC	0.0911	50.85%	0.0902	51.90%
AAPL	0.0223	55.84%	0.0232	68.42%
ATVI	0.0423	54.49%	0.0439	60.53%
SBUX	0.0231	40.49%	0.0213	50.27%

As can be discerned from the table, in the case of stocks, the exclusion of sentiment and attention variables from the model leads even to improved forecast performance, particularly in terms of MDA. This finding corroborates the notion that sentiments extracted from social media and search engines may primarily affect relatively new and less capitalized currencies, and may not have an impact on the forecast performance of actual stocks or highly capitalized cryptocurrencies.

6 Conclusion

We conducted a study using a novel dataset that combines cryptocurrencies returns, social media sentiment indexes, Google trends, and volume. Our objective was to investigate whether these variables can improve log-return forecasts through the use of regularized high-dimensional VAR models compared to benchmark models. We also aimed to determine which data sources and types of sentiment or attention measures are most relevant in terms of Granger-causality in High-Dimensional VARs.

Our results indicate that, for a thirty-day forecast, our regularized high-dimensional VAR models perform similarly in terms of Mean Directional Accuracy (MDA) compared to the most popular benchmark models. Furthermore, we find that sentiment and attention measures significantly contribute to achieving better MDA performance for middle and low capitalized cryptocurrencies, but have no impact on highly capitalized cryptocurrencies and stocks (which we used as a robustness check). This finding may be crucial for portfolio selection strategies and investment analysis.

We also investigated Granger causality between the variables and found that the contribution of sentiment and search engine variables is not translated into Granger causality. In most cases, we cannot reject the null hypothesis of (non-)Granger causality. However, we did find an interesting connection between volume and cryptocurrency returns.

References

- N. Aslanidis, A. F. Bariviera, and Ó. G. López. The link between cryptocurrencies and google trends attention. *Finance Research Letters*, page 102654, 2022.
- R. Auer, C. Monnet, and H. S. Shin. Permissioned distributed ledgers and the governance of money. *Available at SSRN 3770075*, 2021.
- M. Balcilar, E. Bouri, R. Gupta, and D. Roubaud. Can volume predict bitcoin returns and volatility? a quantiles-based approach. *Economic Modelling*, 64:74–81, 2017.
- M. Bańbura, D. Giannone, and L. Reichlin. Large bayesian vector auto regressions. *Journal of applied Econometrics*, 25(1):71–92, 2010.
- J. Baumgartner, S. Zannettou, B. Keegan, M. Squire, and J. Blackburn. The pushshift reddit dataset. In *Proceedings of the international AAAI conference on web and social media*, volume 14, pages 830–839, 2020.
- A. Belloni, V. Chernozhukov, and C. Hansen. High-dimensional methods and inference on structural and treatment effects. *Journal of Economic Perspectives*, 28(2):29–50, 2014.
- B. S. Bernanke, J. Boivin, and P. Eliasch. Measuring the effects of monetary policy: a factor-augmented vector autoregressive (favar) approach. *The Quarterly journal of economics*, 120(1):387–422, 2005.
- C. Catalini and J. S. Gans. Some simple economics of the blockchain. *Communications of the ACM*, 63(7):80–90, 2020.
- L. Catania and S. Grassi. Forecasting cryptocurrency volatility. *International Journal of Forecasting*, 2021.
- L. Catania, S. Grassi, and F. Ravazzolo. Forecasting cryptocurrencies under model and parameter instability. *International Journal of Forecasting*, 35(2):485–501, 2019.
- S. Chan, J. Chu, Y. Zhang, and S. Nadarajah. An extreme value analysis of the tail relationships between returns and volumes for high frequency cryptocurrencies. *Research in International Business and Finance*, 59:101541, 2022.
- P. Ciaian, M. Rajcaniova, and d. Kancs. The economics of bitcoin price formation. *Applied economics*, 48(19):1799–1815, 2016.
- T. Clark and M. McCracken. Advances in forecast evaluation. *Handbook of economic forecasting*, 2:1107–1201, 2013.
- S. Corbet, B. Lucey, A. Urquhart, and L. Yarovaya. Cryptocurrencies as a financial asset: A systematic analysis. *International Review of Financial Analysis*, 62:182–199, 2019.
- K. Daniel, D. Hirshleifer, and S. H. Teoh. Investor psychology in capital markets: Evidence and policy implications. *Journal of monetary economics*, 49(1):139–209, 2002.

- F. X. Diebold and R. S. Mariano. Comparing predictive accuracy. *Journal of Business & economic statistics*, 20(1):134–144, 2002.
- E. F. Fama. Efficient capital markets: A review of theory and empirical work. *The journal of Finance*, 25(2):383–417, 1970.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.
- N. Gandal and H. Halaburda. Can we predict the winner in a market with network effects? competition in cryptocurrency market. *Games*, 7(3):16, 2016.
- M. Glenski, T. Weninger, and S. Volkova. Improved forecasting of cryptocurrency price using social signals. *arXiv preprint arXiv:1907.00558*, 2019.
- C. W. Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*, 37:424–438, 1969.
- C. W. Granger. Testing for causality: a personal viewpoint. *Journal of Economic Dynamics and control*, 2:329–352, 1980.
- H. Halaburda, G. Haeringer, J. S. Gans, and N. Gandal. The microeconomics of cryptocurrencies. Technical report, National Bureau of Economic Research, 2020.
- D. Harvey, S. Leybourne, and P. Newbold. Testing the equality of prediction mean squared errors. *International Journal of forecasting*, 13(2):281–291, 1997.
- A. Hecq, L. Margaritella, and S. Smeekes. Granger causality testing in high-dimensional vars: a post-double-selection procedure. *arXiv preprint arXiv:1902.10991*, 2019.
- N. A. Hitam and A. R. Ismail. Comparative performance of machine learning algorithms for cryptocurrency forecasting. *Ind. J. Electr. Eng. Comput. Sci*, 11(3):1121–1128, 2018.
- C. Hutto and E. Gilbert. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the international AAAI conference on web and social media*, volume 8, pages 216–225, 2014.
- P. C. Jain and G.-H. Joh. The dependence between hourly prices and trading volume. *Journal of Financial and Quantitative Analysis*, 23(3):269–283, 1988.
- S.-J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky. An interior-point method for large-scale ℓ_1 -regularized least squares. *IEEE journal of selected topics in signal processing*, 1(4):606–617, 2007.
- O. Kraaijeveld and J. De Smedt. The predictive power of public twitter sentiment for forecasting cryptocurrency prices. *Journal of International Financial Markets, Institutions and Money*, 65:101188, 2020.
- G. Llorente, R. Michaely, G. Saar, and J. Wang. Dynamic volume-return relation of individual stocks. *The Review of financial studies*, 15(4):1005–1047, 2002.

- D. Miller and J.-M. Kim. Univariate and multivariate machine learning forecasting models on the price returns of cryptocurrencies. *Journal of Risk and Financial Management*, 14(10):486, 2021.
- M. Naeem, K. Saleem, S. Ahmed, N. Muhammad, and F. Mustafa. Extreme return-volume relationship in cryptocurrencies: Tail dependence analysis. *Cogent Economics & Finance*, 8(1):1834175, 2020.
- S. Nakamoto. Bitcoin: A peer-to-peer electronic cash system. *Decentralized Business Review*, page 21260, 2008.
- M. A. Nasir, T. L. D. Huynh, S. P. Nguyen, and D. Duong. Forecasting cryptocurrency returns and volume using search engines. *Financial Innovation*, 5(1):1–13, 2019.
- X. Sun, M. Liu, and Z. Sima. A novel cryptocurrency price trend forecasting model based on lightgbm. *Finance Research Letters*, 32:101084, 2020.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- K. D. West. Asymptotic inference about predictive ability. *Econometrica: Journal of the Econometric Society*, pages 1067–1084, 1996.
- W. Zhang, P. Wang, X. Li, and D. Shen. Multifractal detrended cross-correlation analysis of the return-volume relationship of bitcoin market. *Complexity*, 2018, 2018.

A Appendix

Table 9: Google Trends collected from January 2018 to January 2022.

Google Trends	
<i>binance coin</i>	<i>zen coin</i>
<i>crypto crash</i>	<i>coinbase</i>
<i>miner</i>	<i>BTC</i>
<i>USD coin</i>	<i>hard fork</i>
<i>cardano</i>	<i>stablecoin</i>
<i>cryptocurrency</i>	<i>cold storage</i>
<i>minted</i>	<i>coin</i>
<i>solana</i>	<i>hash</i>
<i>feathercoin</i>	<i>tether</i>
<i>cryptography</i>	<i>litecoin</i>
<i>dogecoin</i>	<i>enjin coin</i>
<i>digital gold</i>	<i>hashing</i>
<i>Bitcoin</i>	<i>token</i>
<i>NFT</i>	<i>crypto</i>
<i>ripple</i>	<i>ICO</i>
<i>private key</i>	<i>coin market</i>
<i>peercoin</i>	<i>soft fork</i>
<i>namecoin</i>	<i>block producer</i>
<i>satoshi</i>	<i>distributed ledger</i>
<i>hot wallet</i>	<i>digital fiat</i>
<i>blockchain</i>	<i>consensus</i>
<i>ethereum</i>	

Table 10: The table shows the cross-correlation coefficients between Twitter and Reddit sentiment as $x_{t\pm h}$ with $h = 0, \pm 1, \dots, \pm 14$ and respective Cryptocurrencies as y_t . All variables used are stationary and scaled and cross-correlation is computed over the whole sample. Bold numbers represent significant correlations. * significance at the 10% level; ** significance at the 5% level; *** significance at the 1% level.

	binance.twitter	bitcoin.twitter	ejun.twitter	ethereum.twitter	feathercoin.twitter	horizon.twitter	litecoin.twitter	manecoin.twitter	peercoin.twitter	terceron.twitter	teletwitter	binance.reddit	bitcoin.reddit	ejun.reddit	ethereum.reddit	feathercoin.reddit	horizon.reddit	litecoin.reddit	manecoin.reddit	peercoin.reddit	terceron.reddit	feather.reddit
-14	-0.018	0.067	0.041	-0.015	-0.056	0.021	-0.045*	0.028	0.003	0.026	0.026	0.005	0.006	-0.002	0.003	0	-0.005	0.004	0.003	0.028	0.028	-0.033
-13	-0.001	-0.045*	-0.016	-0.022	0.021	-0.026	0.046*	-0.014	0.025	0.025	0.025	-0.004	-0.015	-0.017	-0.011	0.005	0.002	-0.018	0.004	0.004	-0.012	0.032
-12	0.033	-0.006	-0.009	-0.033	0.023	-0.027	-0.006	-0.009	-0.021	-0.053**	0.029	0.029	0.022	0.03	0.034	-0.031	-0.029	0.026	0.023	0.043*	0.043*	-0.035
-11	-0.005	-0.033	-0.035	-0.036	-0.032	0.03	-0.045*	0.026	0.011	0.025	-0.039	-0.049*	-0.025	-0.039	-0.021	0.017	0.024	-0.004	-0.02	-0.041	0.005	-0.026
-10	-0.022	0.026	0.002	0.055**	0.015	-0.003	-0.005	0.016	-0.034	0.002	0.005	0.005	-0.019	-0.014	-0.007	0.047*	-0.037	0.001	0.028	-0.016	-0.026	-0.026
-9	0.009	0.023	0.027	-0.016	0.006	0.001	0.046*	-0.008	0.051**	0.033	0.033	0.033	0.033	0.033	0.047*	-0.001**	0.067***	0.004	-0.042	0.011	0.085***	-0.056**
-8	0.006	-0.017	-0.017	0.019	0.009	-0.028	-0.028	-0.028	-0.026	-0.019	-0.019	0.001	0.035	-0.016	-0.015	-0.019	-0.007	-0.016	0.015	-0.012	0.054**	-0.009
-7	-0.009	-0.001	0.001	-0.063**	-0.046*	0.019	0.089***	-0.026	-0.012	-0.059**	0.016	0.016	-0.028	-0.002**	-0.055**	0.044*	-0.022	-0.025	-0.021	0.054**	0.054**	-0.009
-6	-0.027	0.002	-0.01	-0.004	0.025	-0.096	-0.026	0.025	-0.021	0.031	0.038	0.037	0.037	0.029	0.019	-0.005**	-0.02	0.033	0.018	-0.006	0.042	-0.045*
-5	0.017	0.007	0.055**	0.009	0.039	0.016	-0.051**	-0.019	0.043*	-0.033	-0.034	-0.034	-0.016	0.005	0.025	0.058**	-0.003	-0.036	0.051**	-0.026	0.006	0.005
-4	0.014	0.018	-0.02	0.005	-0.021	0.002	0.033	-0.024	-0.028	0.045*	-0.025	-0.025	-0.005	0.005	0.006	-0.029	0.036	0.021	-0.083***	-0.022	0.003	0.003
-3	-0.02	-0.011	0.007	0.051**	0.015	-0.023	0.01	0.029	-0.031	0.004	0.027	0.01	-0.009	-0.009	-0.01	0.036	-0.055**	0	0.045*	-0.023	0.005	0.005
-2	0.011	0.002	0.02	-0.044*	-0.001	0.027	-0.058**	0.01	0.039	-0.018	0.002	0.059**	0.004	0.004	0.029	-0.049*	0.048*	0.012	0.002	0.077***	0.006	0.006
-1	-0.033	-0.049*	-0.009	-0.009	-0.019	0.058**	0.022	-0.026	-0.01	0.054**	0.001	-0.03	0.019	-0.042	-0.042	0.046*	-0.004	0.014	-0.022	0.001	0.035	-0.037
0	0.019	0.029	0.038	0.036	0.023	-0.037	0.026	0.036	0.03	-0.067***	0.011	0.007	0.087***	-0.002	0.028	0.021	0.027	0.031	0.03	-0.053**	-0.037	-0.037
1	0.003	0.024	0.003	0.014	-0.027	-0.002	0.013	-0.025	0.024	0.031	-0.007	0.018	0	0.046*	-0.029	-0.026	-0.014	-0.026	-0.003	0.007	0.007	0.007
2	-0.021	-0.027	-0.019	-0.009	0.026	-0.022	-0.017	0.004	-0.017	-0.008	-0.026	-0.026	-0.012	-0.011	-0.04	0.036	-0.027	-0.023	0.038	0.04	0	0
3	0.027	0.01	0.018	0.01	0.022	0.013	-0.002	-0.01	-0.039	-0.028	-0.028	0.003	0.003	-0.001	0.006	-0.015	0.015	-0.002	-0.015	0.005	-0.006	-0.006
4	-0.029	-0.051**	-0.01	-0.027	-0.022	-0.032	0.009	0.043*	0.041	0.018	0.045*	0.011	0.006	0.006	0.002	0.034	-0.012	-0.01	-0.026	-0.049*	-0.004	-0.004
5	0.038	0.026	0.006	-0.008	-0.012	0.033	-0.041	-0.076***	-0.011	0.016	0.002	-0.051**	-0.012	-0.016	-0.001	-0.048*	0.006	0.016	0.03	0.032	0.012	0.012
6	-0.031	0.013	-0.03	0.001	0.016	0.024	0.018	0.060**	0.021	0.02	0.009	0.009	-0.025	0.027	-0.005	0	0.002	0.005	-0.014	0.004	0.004	-0.012
7	0.024	0.02	0.021	0.038	0.012	0.006	0.032	-0.027	-0.005	-0.027	0.033	0.033	-0.02	-0.029	-0.03	0.014	0.001	0.012	0.02	0.021	-0.017	-0.017
8	-0.044*	-0.052**	0.014	-0.034	-0.080***	0.014	-0.015	-0.011	-0.053**	-0.059**	0.009	-0.01	0.035	0	-0.019	0.006	-0.053**	-0.017	-0.004	0.044*	-0.017	-0.017
9	0.034	0.018	-0.026	0.017	0.056**	-0.013	0.016	0.016	0.022	0.066**	-0.02	-0.001	-0.011	-0.011	0.007	0.033	-0.017	0.050*	0.004	0.014	0.014	-0.012
10	-0.01	0.03	0.02	-0.019	-0.006	-0.028	-0.027	0.01	-0.003	0.028	-0.003	0.022	0.008	0.008	-0.034	-0.019	0.005	-0.013	-0.015	-0.01	0.005	0.005
11	0.041	-0.057**	-0.006	0.018	0	0.018	-0.017	-0.016	0.048*	-0.061**	-0.02	0.007	0.007	0.007	0.006	-0.008	-0.008	-0.023	-0.001	-0.019	-0.029	-0.029
12	-0.053**	0.014	-0.015	-0.032	0.024	-0.044*	-0.023	0.053**	-0.018	0.082***	0.012	-0.031	-0.009	-0.009	-0.01	-0.022	0.028	-0.02	-0.009	-0.005	0.019	0.019
13	0.002	-0.014	0.012	0.062**	-0.018	0.025	0.052**	-0.036	-0.002	-0.060**	0.043*	-0.011	-0.029	-0.022	0.046*	-0.027	-0.007	-0.007	0.021	-0.014	-0.017	-0.017
14	0.024	0.034	0.007	-0.021	-0.019	0.009	-0.011	-0.031	0.041	0.039	-0.018	0.046	0.019	0.007	-0.052**	0.018	0.021	-0.001	0.04	0.04	0.04	0.04

