



**Faculty of
Mathematics
and Informatics**

VILNIUS UNIVERSITY
FACULTY OF MATHEMATICS AND INFORMATICS
MASTER'S STUDY PROGRAMME
MODELLING AND DATA ANALYSIS

Evaluation of Consumer Confidence Indicators Using Social Media and Administrative Data

Master's thesis

Author: Akvilė Vitkauskaitė

VU email address: akvile.vitkauskaite@mif.stud.vu.lt

Supervisor: Andrius Čiginas

Vilnius

2024

Abstract

The main objective of this study is to nowcast and forecast the Consumer Confidence Index (CCI). The aim is to estimate the current month's CCI values faster than those obtained using the traditional survey methodology, which usually provides results at the end of the month. For instance, while the official CCI for November would typically be available in the last few days of November, this research aims to provide an early estimate at the beginning of November, utilizing data collected at the start of the month. This is achieved by combining key economic indicators with historical CCI values. The research includes examining the relationship between traditional survey-based indicators and consumer sentiment expressed on social media platforms. Social media expressions, particularly from X (Twitter), are analyzed through its official API. The sentiment analysis of tweets has enabled us to create a Social Media Indicator (SMI) that offers a distinct advantage in our predictive models. In addition, the study explores the possibility of integrating key economic indicators from administrative data, such as inflation rate, income statistics, and unemployment. In general, obtaining data for research from popular social platforms such as Facebook and Instagram is not possible due to stringent privacy policies and data protection regulations. Nevertheless, data are easily and legally available from X, but this platform is not so popular in Lithuania. Therefore, the representativeness of X data raises special issues. Taking everything into account, by combining traditional economic indicators with advanced sentiment analysis from X, the study seeks to deliver prompt CCI predictions ahead of standard survey timelines.

Keywords: Consumer confidence, Social Media, Twitter, Sentiment Analysis, SARIMAX, VECM, Random Forest, XGBoost

Santrauka

Šio tyrimo pagrindinis tikslas — prognozuoti vartotojų pasitikėjimo rodiklį (VPR). Siekiama įvertinti einamojo mėnesio VPR reikšmes greičiau nei tai daroma naudojant tradicinius apklausų metodus, kurie paprastai rezultatus pateikia mėnesio pabaigoje. Pavyzdžiui, nors oficialus lapkričio mėnesio VPR paprastai būna skelbiamas lapkričio mėnesio paskutinėmis dienomis, šis tyrimas siekia pateikti ankstyvą įvertį jau lapkričio pradžioje, naudojant duomenis, surinktus mėnesio pirmomis dienomis. Tai pasiekama derinant pagrindinius ekonominius rodiklius su istorinėmis VPR reikšmėmis. Tyrimas apima ir sąryšių tarp tradicinių, apklausomis grindžiamų rodiklių bei vartotojų nuomonės, išreikštos socialiniuose tinkluose, analizę. Vartotojų sentimentai yra gaunami iš socialinio tinklo X (Twitter) naudojant oficialių API. Twtiterio žinučių sentimentų analizė leidžia mums sukurti socialinės medijos indikatorius (SMI), kuris yra naudingas prognozavimo modeliams, nes padidina prognozės tikslumą, atspindėdamas socialiniuose tinkluose vyraujančias nuotaikas. Šiame darbe iš administracinių duomenų bandome integruoti pagrindinius ekonominius rodiklius: infliacijos lygį, pajamų statistiką ir nedarbo apimtį. Verta paminėti, kad gauti duomenis iš populiarių socialinių platformų, tokių kaip Facebook ir Instagram, neįmanoma dėl griežtos privatumo politikos ir duomenų apsaugos reglamentų. Tuo tarpu, duomenys iš X yra lengvai ir legaliai prieinami, tačiau Lietuvoje ši platforma nėra tokia populiari. Todėl kyla tam tikrų X duomenų reprezentatyvumo problemų. Taigi, šis tyrimas siekia pateikti VPR prognozes anksčiau nei rezultatai gaunami tradiciniais metodais, derinant tradicinius ekonominius rodiklius ir pažangią sentimentų analizę X duomenims.

Raktiniai žodžiai: Vartotojų pasitikėjimo rodiklis, Socialiniai tinklai, Twitter, Sentimentų analizė, SARIMAX, VECM, Random Forest, XGBoost

1 Introduction

Consumer confidence is a vital economic indicator that influences the decision-making processes of policymakers, businesses, and investors, providing valuable insights into individuals' sentiments and expectations regarding the state of the economy [21]. Traditionally, consumer confidence indicators have been derived from survey data, capturing the opinions and perceptions of individuals through structured questionnaires [31, 12, 36].

With the rise of social media platforms and the abundance of user-generated content, there is a new opportunity to assess consumer sentiment in new ways. Platforms like internet blogs, discussion forums, and social media sites enable individuals to openly express their opinions and experiences [1]. By analyzing this rich data through sentiment analysis techniques, valuable insights can be gained regarding public opinion and attitudes, including the evaluation of consumer confidence. These insights can inform decision-making in various fields such as market research and business strategies.

In addition to social media data, administrative data from government agencies, such as employment statistics, income records, and other relevant economic indicators, can provide a rich and objective source of information. By incorporating these auxiliary variables, it can potentially improve the accuracy and timeliness of consumer confidence indicators. The study by Curtin (2007) [11] emphasizes the significant relationship between changes in the unemployment rate and consumer sentiment, highlighting the powerful influence of employment conditions on consumer expectations.

This master thesis aims to evaluate the effectiveness and reliability of consumer confidence indicators derived using social media data and administrative data. Also, an examination of the relationships between consumer sentiment expressed on social media platforms and traditional survey-based indicators is performed. Furthermore, the potential of using administrative data as auxiliary variables to enhance the forecasting accuracy of consumer confidence indicators is explored.

We aim to bridge the gap between traditional survey-based indicators and emerging data sources by evaluating and utilizing social media and administrative data to measure and forecast consumer confidence. The objective is to provide valuable insights and contribute to the advancement of economic forecasting and decision-making in an increasingly digitalized world.

Our study used SARIMAX, VECM, Random Forest, and XGBoost to predict the CCI from social media and administrative data. XGBoost yielded the highest accuracy, SARIMAX followed closely, and Random Forest and VECM were comparatively less precise. We are going to present a detailed analysis and comparison of these models.

The thesis is divided into several sections. Section 2 provides a literature review that examines the historical context and evolution of consumer confidence indicators, the emerging role of social media data, and the inclusion of administrative data as auxiliary sources. Section 3 introduces the data used in this study, detailing the sources, collection methods, and pre-processing steps. Section 4 conducts an exploratory analysis. Section 5 describes the evaluation metrics and rolling forecasting method, setting the stage for Section 6, which discusses the development and implementation of various forecasting models. Finally, Section 7 summarizes the findings, compares the models' performance, and suggests directions for future research. Appendices include supplementary figures, tables, and Python code.

2 Literature review

2.1 Historical development of consumer confidence indicators

In the middle of the twentieth century, consumer confidence measurements became highly important for predicting economic trends. They gained recognition among businesses and government economists as a reliable method to forecast market behaviour and understand consumer choices. This led to a greater emphasis on comprehending consumers, including their attitudes, expectations, and the psychological factors that shape their economic decisions [28].

The development of the CCI requires the collaboration of various organizations and institutions. George Katona and Rensis Likert pioneered the initial methods to gauge consumer confidence in the late 1940s, aiming to integrate tangible measurements of expectations into models analyzing spending and saving patterns [23]. Mueller (1963) [31] played a significant role in advancing consumer confidence indices by introducing the concept of using them to forecast consumption patterns. Mueller's innovative research involved analyzing data from the Michigan survey of consumers over a period of ten years. Her findings demonstrated the vital role of consumer confidence in explaining spending habits, as evidenced by its inclusion in a regression model that considered previous consumption levels [31].

The Consumer Confidence Board in the USA has played a crucial role in the development of consumer confidence indices, with one notable index being the CCI introduced in 1967. Their survey evaluates individuals' perspectives on current and future economic conditions, as well as their employment prospects. The CCI is widely recognized as an indicator of the strength and stability of the U.S. economy. In May 2021, the Conference Board Consumer Confidence Survey transitioned from being administered by The Nielsen Company to Toluna, a technology company with a large consumer panel of over 36 million individuals. Before November 2010, the survey was conducted by TNS through mail [36].

The contemporary administration of the CCI in the EU lies under the purview of Eurostat, the statistical office of the European Union, with monthly calculations conducted across all member countries [16]. Since May 2001, the State Data Agency (Statistics Lithuania) [26] has been responsible for conducting similar surveys in Lithuania. The main objective of the consumer opinion statistical survey is to obtain information regarding consumers' intentions to make purchases, their saving capabilities, as well as their perceptions of the economic situation and its influence on their intentions.

2.2 The role of social media data in estimating consumer confidence

The integration of social media data in the estimation of consumer confidence indicators has gained significant attention in recent years. Innovative and accelerated use of big data sources has become particularly important in addressing challenges, exemplified by the increased frequency and application of Twitter data in official statistics during the COVID-19 crisis [5]. Alternative data sources, including social media platforms such as Twitter and Facebook, have emerged as valuable resources for capturing real-time information reflecting public sentiment and opinions [18, 37]. In the realm of official statistics, notable efforts have been made to explore the use of social media data for constructing consumer confidence indices. Istat's Social Mood on Economy Index in Italy, developed by Catanese et

al (2022)[8] incorporates social media data alongside traditional survey-based measures. Analyzing the sentiment expressed in tweets related to the economy, this index provides a complementary perspective on consumer confidence. Additionally, van den Brakel et al. (2017) [6] investigated the utilization of alternative data sources in the Dutch Consumer Confidence Survey. Their research developed a multivariate structural time series modelling methodology that incorporates social media data and repeated survey data. This approach aims to enhance the accuracy and timeliness of consumer confidence estimates, even in the absence of sample data.

While there is potential for using social media data in estimating consumer confidence in official statistics, challenges such as data noise, representativeness, and the development of reliable methodologies still need to be addressed [35]. Ongoing research efforts [4] continue to explore innovative approaches, proposing a novel method of social media analytics called NACOP. By utilizing big data and data science techniques, NACOP analyzes social media data on purchasing behaviour, jobs/employment, consumer price, and personal finance. This approach aims to provide more accurate insights into the true state of the economy, offering significant implications for researchers and practitioners.

2.3 Administrative data as auxiliary data for estimating CCI

Nowzohour and Stracca (2017) [33] find that the CCI is highly correlated with expectations on the general economic situation, unemployment, and the financial situation. Consumer confidence is positively associated with future inflation, industrial production growth, real house price growth, and real appreciation, while negatively associated with unemployment and short-term interest rates. Demirel and Artan (2017) [14] employ panel causality analysis to examine the causality relationships between economic confidence and fundamental macroeconomic indicators in selected European Union countries. A unidirectional causality is observed from real exchange rate and interest rate to economic confidence, and from economic confidence to unemployment. The results highlight the significant impact of confidence on production, consumption, inflation, and unemployment. The findings emphasize the importance of confidence as a psychological factor in economic decision-making.

The previously mentioned research provides insights and inspiration to explore the use of auxiliary data, such as unemployment data, as a means to enhance the accuracy of CCI forecasting models. The observed correlations between these auxiliary variables and the CCI suggest the potential benefits of incorporating them into the forecasting process.

3 Data and dataset

3.1 Overview

The dataset used in this study is designed to forecast the CCI faster than traditional approaches. CCI (Fig. 1) is assessed by a survey and it is calculated as the arithmetic mean of four balanced estimates: changes in household financial situations over the past and upcoming 12 months, anticipated changes in the national economy over the next 12 months, and variations in spending on significant purchases (like furniture or household appliances) compared to the previous year [26]. This dataset combines public sentiment from social media with key economic indicators. It contains monthly data from 2018 until November 2023, including the COVID-19 pandemic and geopolitical developments - Russia’s invasion of Ukraine in February 2022. The observed declines in the CCI correspond to these periods of global and regional turbulence.

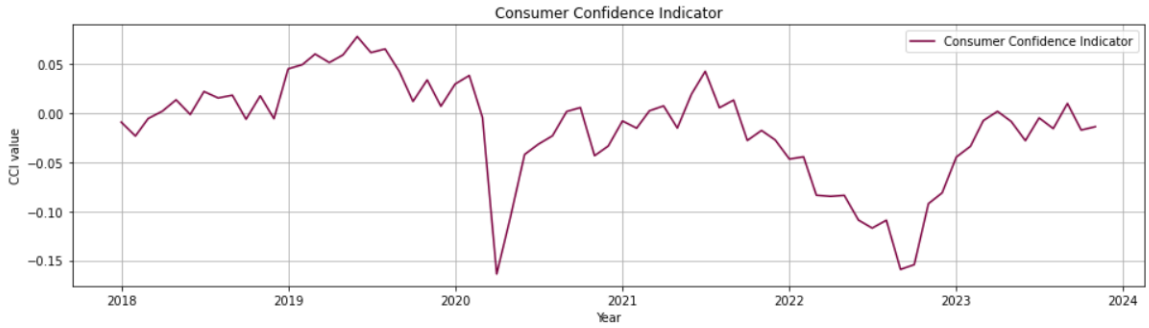


Figure 1: Consumer confidence indicator

3.2 Social media data (SMI)

We aim to refine the nowcasting of Lithuania’s CCI by integrating social media data. Central to this integration is the Social Media Indicator (SMI), an index constructed from X data that quantifies public economic sentiment, offering real-time insights to complement traditional CCI measures. The social media component, derived from X using the platform’s official API, captures the economic sentiment of Lithuanians through a crafted set of economic keywords (see Table 4).

As an auxiliary variable, the SMI enriches the dataset to increase the accuracy of the CCI forecasts.

3.2.1 Data gathering and preprocessing

In the middle of 2023, tweets reflecting Lithuania’s economic sentiment began to be collected. The comprehensive list of 101 economic keywords and their variations was created from the words of official survey questionnaires and related research. Using this list, tweets were gathered weekly. To overcome the API’s limitation of retrieving only the most recent 7-day-old tweets based on these keywords, the timelines method was used. This approach involves extracting the historical tweets directly from the timelines of users who had recently tweeted with our selected keywords, based on the assumption that their past tweets would likely contain similar economic content.

Before sentiment analysis is applied, the collected tweets go through several preprocessing steps. These steps are critical to ensure data quality and suitability for sentiment analysis. Initially, Twitter handlers, hashtags, URLs, special characters, and single characters are removed to clean the text. After that, multiple spaces are replaced with a single space, and all empty lines are replaced with NaNs (special floating point values in Python that are used to represent undefined or unrepresentable values). After dropping rows with missing values (NaNs), we excluded non-Lithuanian tweets (as Lithuanian and other languages have linguistic similarities, some non-Lithuanian tweets were collected).

3.2.2 Sentiment analysis

The final but not least step in developing the SMI is the sentiment analysis. Sentiment analysis, known as opinion mining, is a part of natural language processing (NLP). It identifies and categorises the opinions expressed in a text, particularly to determine whether the writer's view on a particular topic is positive, negative, or neutral. The purpose of sentiment analysis is to understand the attitudes, emotions and opinions from textual data [27].

Each tweet's text was analysed using five different techniques: TextBlob, Vader, AFINN, Flair, and Transformers. These analyzers evaluate the emotional sentiment of the text and assign scores that classify the sentiment as positive, neutral, or negative, except for Flair, which distinguishes only between positive and negative. Using multiple methods provides a broader and more nuanced perspective on public sentiment. Each different analyzer's sentiment index has slightly different relationships with the CCI.

The average sentiment score is calculated for Twitter users who tweet more than once a month. This method treats each tweet equally because each tweet has the same weight. This provides a balanced view of that user's economic sentiment throughout the month, no matter how often they tweet.

Each user's monthly sentiment scores are then categorized as either positive, neutral or negative (for Flair: positive or negative). The final SMI for the month is determined by calculating the balance of these sentiments. The number of negative expressions is subtracted from the total number of positive expressions. This figure is then divided by the sum of all positive and negative expressions, converting it into a percentage. Then SMI is adjusted to match the scale of the CCI, as the raw SMI values are typically ten times larger than the CCI values.

The interpretation of the Social Media Indicator (SMI) is that a positive sentiment score indicates an optimistic public mood regarding the economy, while 0 reflects a neutral mood, a negative - a pessimistic public sentiment.

This calculation is performed for each month in the dataset, spanning from January 2018 to November 2023 (71 months). As a result, five separate SMI indices emerge, each corresponding to one of the sentiment analysis techniques employed.

3.2.3 Sentiment analysis tools

Here we provide a brief overview of each sentiment analysis tool based on the general methodologies they employ.

- **TextBlob.** TextBlob simplifies text processing in Python, with sentiment analysis that evaluates text polarity and subjectivity. Polarity is a float between $[-1, 1]$, where -1 signifies negative and +1 positive sentiment. Subjectivity, ranging from $[0, 1]$, indicates how much personal opinion or emotion is expressed in the text. It is most suited for general text analysis including social media and product reviews [29].
- **Vader.** Vader (Valence Aware Dictionary for sEntiment Reasoning) is a rule-based sentiment analysis tool that is designed for social media sentiment analysis. It uses a combination of a sentiment lexicon and grammatical rules. It scores texts by evaluating slangs, emojis, and colloquial language, with a compound score indicating the overall sentiment. This compound score is normalized to range between $[-1, 1]$, where higher positive values denote more positive sentiment. It is effective for tweets and online reviews [19].
- **Transformers.** In this study, we utilize the Transformers library, developed by Hugging Face, which uses advanced models like BERT (Bidirectional Encoder Representations from Transformers) for deep contextual understanding. Specifically, we employ the "finiteautomata/bertweet-base-sentiment-analysis" model from the Transformers library for sentiment analysis. This particular model is a specialized adaptation of BERT, explicitly tailored for analyzing Twitter data. It effectively outputs a sentiment classification – positive, negative, or neutral – based on its contextual interpretation of the text, making it highly suited for understanding the unique linguistic characteristics of tweets [39].
- **Flair.** We utilize the Flair 'en-sentiment' model which is a pre-trained model for sentiment analysis provided by the Flair framework. This model excels in interpreting the sentiment of Twitter data. It provides a detailed sentiment score that reflects the intricacies of language use on social media platforms. Flair's approach to sentiment analysis is comprehensive as it gives a sentiment label along with a confidence score to signify the model's certainty. This method is particularly valuable when analyzing texts that require an understanding of deep contextual meanings, such as the brief and dynamic content of tweets [3].
- **Afinn.** Afinn employs a simple lexicon-based approach, where each word in the text is scored with a predefined sentiment value. These scores are then summed up to give an overall sentiment score for the text, with positive scores indicating positive sentiment and negative scores indicating negative sentiment [32]. As it relies on individual word scores without considering context means it does not capture the full context and may miss nuanced or ironic expressions. Despite this, Afinn is highly suitable for real-time sentiment analysis, such as monitoring Twitter during fast-moving events or trends. Its quick processing capability makes it efficient for analyzing straightforward texts and capturing the general sentiment direction, a useful feature for immediate public opinion assessment [10].

3.3 Administrative data

Alongside the SMI, the final dataset is enriched with a set of administrative economic indicators. All data for new indicators development is taken from the State Data Agency's official Database of Indicators while the number of unemployed is taken from the Lithuanian Public Employment Service. These variables are key in providing traditional economic signals. Below are listed indicators of the final dataset with explanations of how each was developed.

- **Average Wage.** It is made from a quarterly indicator, named "Average earnings", which is further disaggregated using linear interpolation into months and expressed in euros. It reflects the gross monthly earnings across the economy, providing insight into the earning power and potential spending capacity of the Lithuanian workforce.
- **Pension Data.** It is developed from the "Average state social insurance old-age pension" which is also a quarterly indicator. Given that significant changes in pension averages tend to occur at the start of the year, the average of the whole quarter was assigned to each month of the quarter possibly with insignificant information loss while preserving the dataset's integrity and having monthly data.
- **Inequality Indicator.** Calculated as the ratio of the average pension to the average wage, this indicator serves as a simple measure of income inequality. It is particularly useful as it potentially stands independent of price level fluctuations.
- **Inflation Rate.** It is the "Consumer Price Index (CPI) - based average annual inflation" indicator when set against the corresponding period of the previous year. This monthly data does not require any additional modifications. It is an extremely important measure that reflects the changing cost of living and its impact on consumer purchasing power.
- **Unemployment without seasonality.** The first unemployment metric included is the "Seasonally adjusted unemployment rate." This figure is reported monthly and reflects the unemployment rate after adjustments for seasonal labour market fluctuations.
- **Unemployment rate.** The second measure of unemployment is calculated by dividing the "number of unemployed at the end of the period" by the "Resident population at the beginning of the month." The numerator is provided by the Lithuanian Public Employment Service, offering a monthly count of individuals actively seeking employment. The denominator, sourced from the State Data Agency, gives the population size at the start of each month, providing a standard base for calculating the unemployment ratio.

These selected variables are combined into the final dataset to provide a comprehensive picture of the economic conditions that influence consumer confidence. The dataset, thus, not only reflects the sentiment captured through social media but is also grounded in the concrete economic realities as represented by these indicators.

3.4 Limitations of data

Twitter's user base in Lithuania, as reported in early 2023, stood at approximately 389,000, accounting for 14.2% of the overall population. This statistic becomes slightly more significant when considering Twitter's user age restriction. With the platform limiting its use to individuals 13 years and older, the 389,000 figure actually represents about 16.4% of the "eligible" audience in Lithuania. These insights, drawn from the "DIGITAL 2023: LITHUANIA" report by Simon Kemp [24], primarily reflect Twitter's advertising reach, which may not fully align with the actual active user base.

Additionally, the inclusion of "non-human" accounts in these figures complicates the true representation of user demographics. The fluctuating nature of Twitter's user metrics, such as the notable increase in 2022 followed by a decrease in early 2023, also necessitates a careful interpretation of these statistics, especially when utilized for economic sentiment analysis.

On the other hand, the integrity and representativeness of the Twitter data, which is critical to this study, was based on a cointegration test with the CCI. Cointegration refers to a statistical relationship where two or more time series, although individually non-stationary (their statistical properties like mean and variance change over time), move together in a long-term, stable manner. The test revealed a meaningful cointegration, confirming that Twitter data is important in assessing user sentiment in Lithuania.

The administrative data, while providing a solid economic backdrop, is not without its limitations. Certain variables like the population and the seasonally adjusted unemployment rate are initially reported as provisional for the last few months and may be subject to later revisions. Additionally, regarding the seasonally adjusted unemployment rate, only the data from the previous month is published at the end of each month. Therefore, for the November data, the same value as that of October is utilized as of the data collection cut-off in the middle of December 2023.

Such limitations, inherent in both social media and administrative data sources, could affect the precision of the nowcasting results. As data sources become more refined and comprehensive (for example, the State Data Agency will provide monthly data along with quarterly data), the potential for achieving a higher degree of accuracy in nowcasting the CCI is anticipated.

4 Explanatory analysis

The heatmap (Fig. 2) provides a visual representation of the strength and direction (correlation) of the relationship between the CCI and a set of variables over different time lags, ranging from 1 to 8 months. The most significant correlation is found with Inflation. It stands out with a strong negative correlation at a short 2-month lag (-0.67), indicating a potentially predictive relationship with the CCI. While sentiment indicators show variable correlations, Transformers_SI and Flair_SI at lag 5, along with Vader_SI at subsequent lags, display promising positive correlations. The Average_wage and Pension correlations suggest that higher wages and pensions may not always correlate with higher consumer confidence, potentially due to the complex nature of economic perception. Both Unemployment_without_seasonality and Unemployment_rate exhibit a negative correlation with CCI, more pronounced in the early lags.

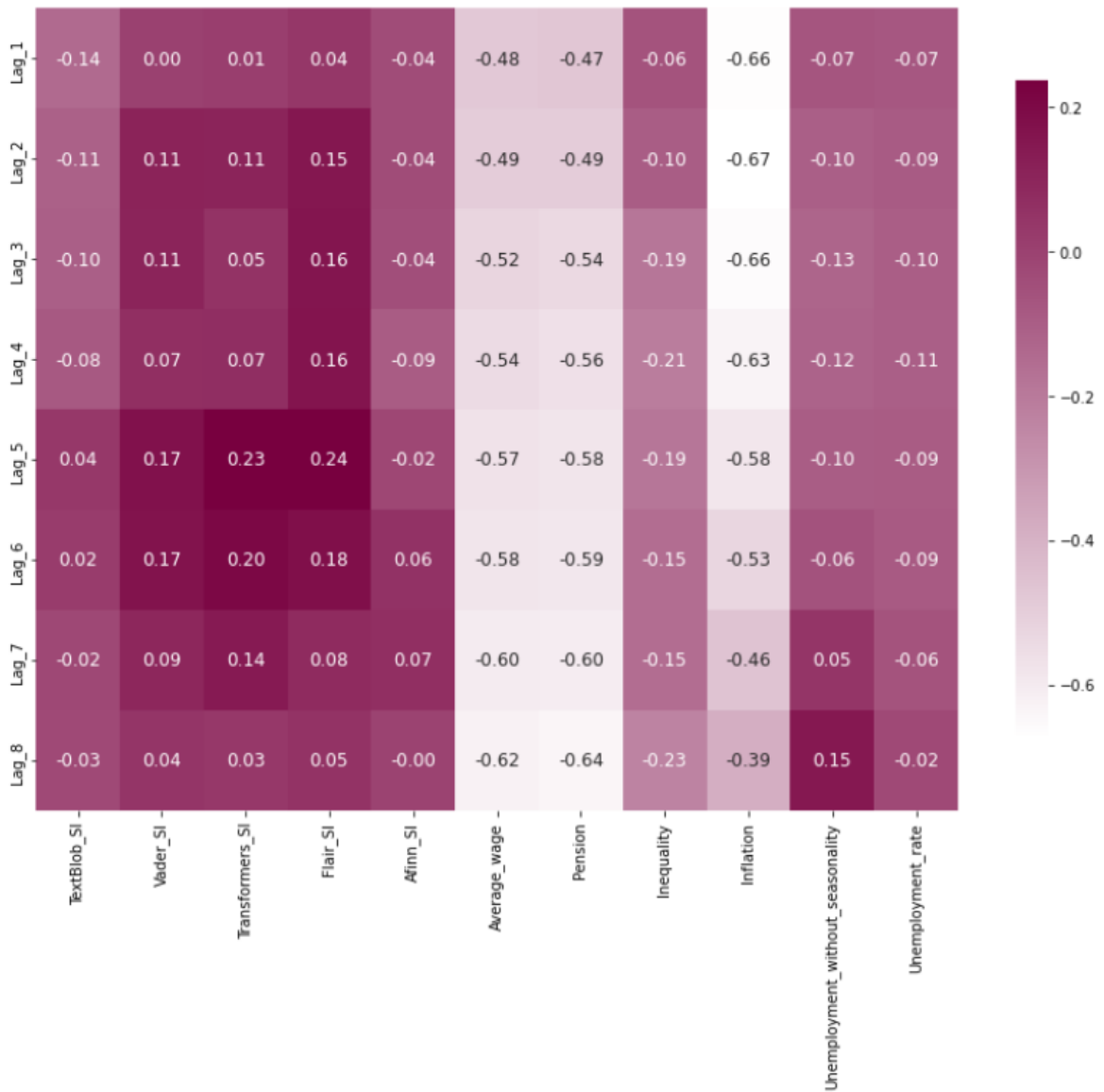


Figure 2: Correlations heatmap of CCI with lagged variables

Fig. 3 presents a visual overview of key economic indicators and sentiment indices. Notably, the CCI reflects a significant dip in April 2020, likely attributed to the initial impact of the COVID-19 pandemic, which officially began in March 2020. This event's economic ramifications are further illustrated by the subsequent peaks in unemployment indicators, which lag the CCI's trough by a few months, suggesting a delayed response in the labour market to the crisis.

Furthermore, following Russia's invasion of Ukraine in February 2022, a downward trend in the CCI can be observed, coinciding with a sharp rise in inflation. This inverse relationship between the CCI and inflation during this period echoes the patterns identified in the correlation heatmap (Fig. 2), underlining the potential impact of geopolitical tensions on consumer confidence and economic stability.

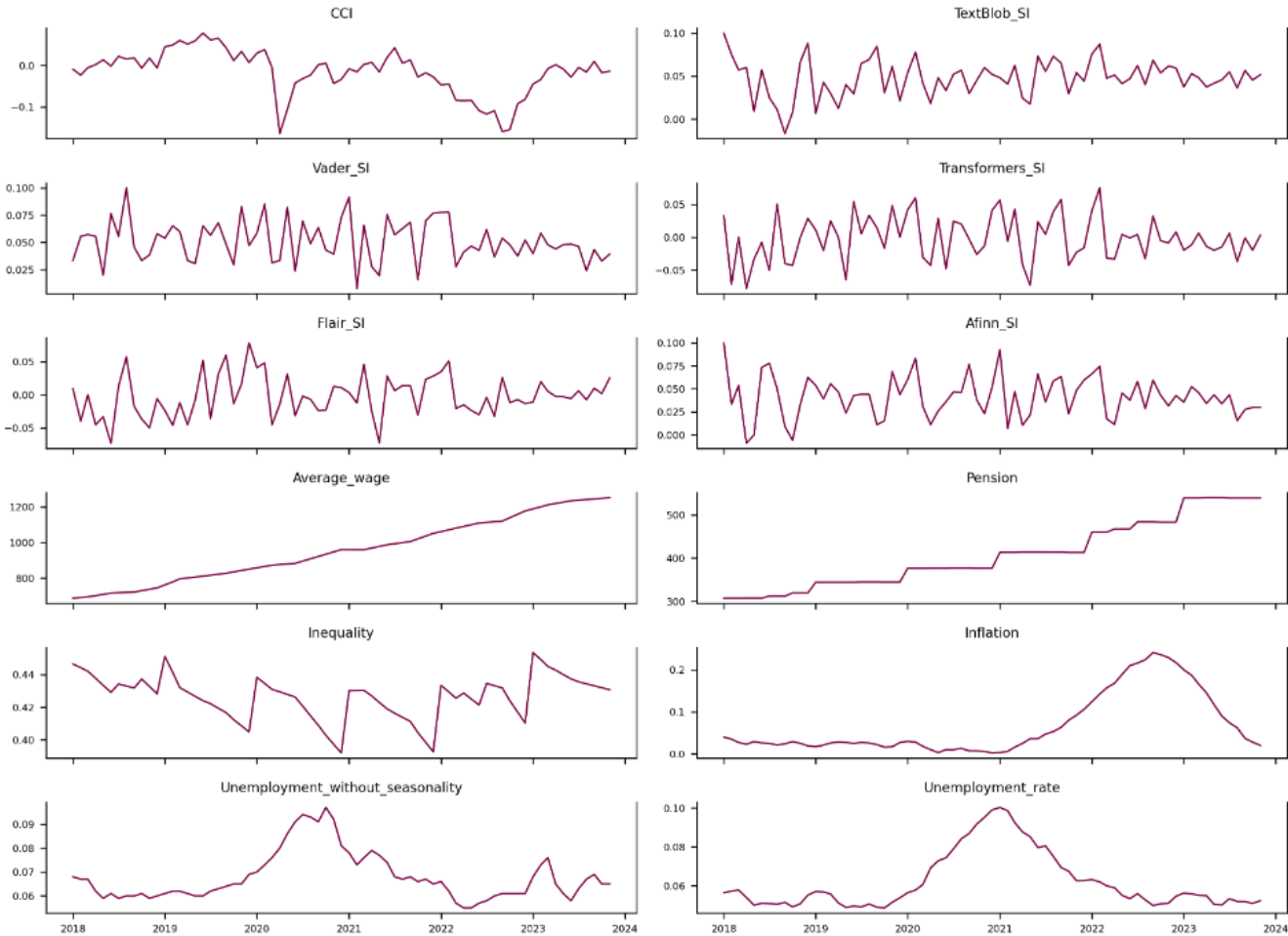


Figure 3: CCI, sentiment indices and key economic indicators

5 Evaluation metrics and rolling forecast approach for forecasting accuracy

To evaluate the accuracy of our forecasting models, we use an effective method called rolling forecast. This method works by predicting the next step and then updating the model with the latest information before making another prediction. It is like making a series of short-term guesses, one after the other, and getting better each time because we use the most recent data.

We evaluate the accuracy of our forecasting models using a set of statistical metrics:

- **Mean Absolute Error (MAE):** MAE measures the average magnitude of errors in predictions, without considering their direction. It's given by the formula:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|, \quad (1)$$

where y_i are the actual values, \hat{y}_i are the predicted values, and n is the number of observations. MAE is a straightforward measure of prediction accuracy.

- **Mean Squared Error (MSE):** MSE calculates the average squared difference between estimated and actual values. The formula for MSE is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2. \quad (2)$$

MSE gauges the variance of forecast errors.

- **Root Mean Squared Error (RMSE):** RMSE is the square root of MSE and measures the standard deviation of prediction errors. Its formula is:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}. \quad (3)$$

RMSE is a scale-dependent accuracy measure often used in forecasting.

- **Akaike Information Criterion (AIC):** AIC is a measure used for model selection that balances the complexity of the model with its fit to the data. The formula for AIC is:

$$\text{AIC} = 2k - 2 \ln(L), \quad (4)$$

where k is the number of parameters in the model and L is the maximum likelihood of the model. AIC is especially valuable for models with numerous parameters, like SARIMAX in our case, as it helps select a model that best explains the data with minimal complexity. A lower AIC value indicates a preferable model balance between accuracy and simplicity [25, 2].

Another one popular measure MAPE (Mean Absolute Percentage Error) is not used as it can be misleading when actual values are near zero, as with CCI values that oscillate around zero. The error metric can be disproportionately inflated, making MAPE unsuitable for our analysis [20].

For each of the forecasting models, we computed the MAE, MSE, and RMSE as key metrics to assess and compare their accuracy in forecasting the CCI. In the case of the SARIMAX model, we additionally calculated the AIC, which is particularly useful for comparing models with different numbers of parameters, providing a balance between model fit and simplicity.

6 Model development and implementation

This section explores 4 models to forecast CCI: SARIMAX, VECM, Random Forest, and XGBoost. Here we explain the reasons why these models were selected and discuss the selection of key exogenous variables and their impact on forecast accuracy. Additionally, we examine the practical application of these models, including their optimization and effectiveness in predicting CCI values. In this section, we aim to provide an understanding of the model development process.

6.1 SARIMAX

6.1.1 SARIMAX model fundamentals

The SARIMAX model, an extension of the well-established ARIMA model, is widely recognized for its effectiveness in time series analysis. As Manigandan et al. (2021) [30] notes, ARIMA is one of the most common methods in this field, and SARIMAX enhances its capabilities by accounting for seasonality and external factors.

The acronym SARIMAX stands for Seasonal AutoRegressive Integrated Moving Average with eXogenous variables. This model is denoted as SARIMAX(p,d,q)(P,D,Q)_s, where 'p' represents the order of the autoregressive (AR) part, 'd' is the order of integration or differencing, and 'q' is the order of the moving average (MA) part. The seasonal components are captured by 'P' (seasonal AR order), 'D' (seasonal differencing order), 'Q' (seasonal MA order), and 's' indicates the frequency or number of observations per seasonal cycle.

A key feature of the SARIMAX model is the inclusion of 'X', signifying exogenous variables. These are external predictors or input variables that can influence the target variable being forecasted. Exogenous variables provide additional external information that can enhance the model's predictive capability [34].

6.1.2 Seasonal and trend decomposition using loess of the CCI

Seasonal and trend decomposition using loess (Locally Estimated Scatterplot Smoothing), in short, STL decomposition, is essential for revealing seasonality in time series, a key factor in applying models like SARIMAX effectively. It separates the data into trend, seasonal, and residual components, allowing us to visually discern seasonal patterns. This method, although not a statistical test, is enough to determine whether a series is seasonal to apply the appropriate model (in this case SARIMAX) for forecasting [34].

STL decomposition (Fig. 4) of the CCI series reveals distinct components that contribute to the overall behaviour of the time series. Observed data shows raw values, while the trend component reflects long-term progress, smoothing out short-term fluctuations. The STL decomposition reveals a noticeable seasonal pattern with noticeable cycles that indicate a strong seasonal influence. Although seasonality is clear, the variation in magnitude and timing of these cycles from year to year suggests that the pattern, while consistent, does not strictly follow a uniform structure.

The residuals representing the remaining variation after accounting for trends and seasonality look randomly distributed.

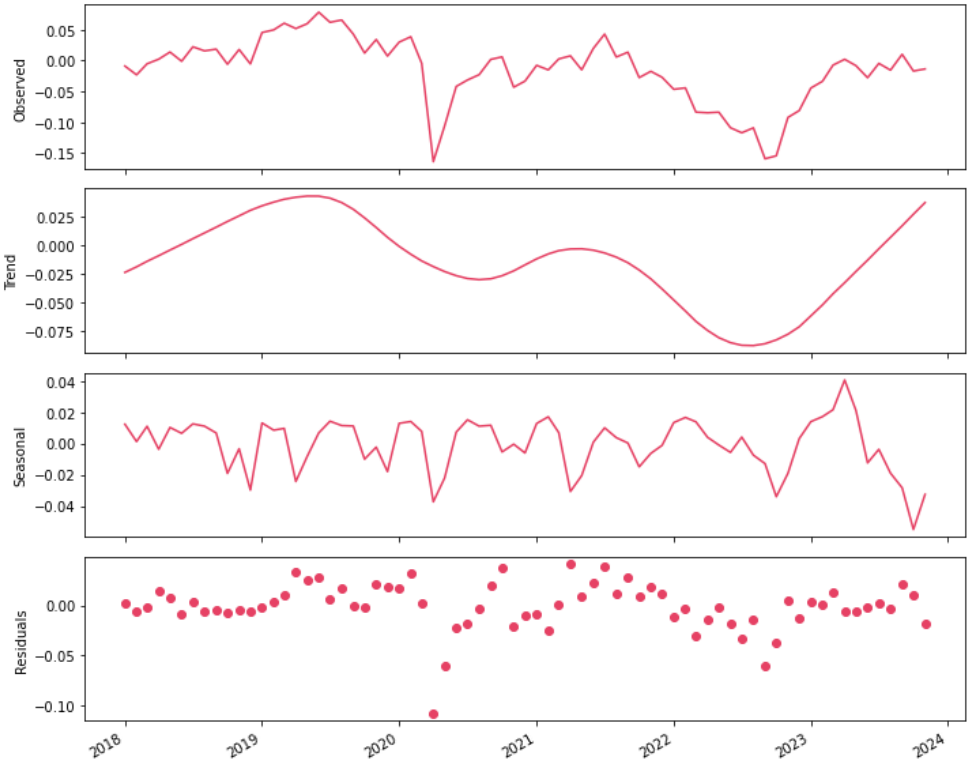


Figure 4: The STL decomposition of the CCI

The decomposition shows some seasonal effects with irregularities, supporting the decision to use the SARIMAX modelling approach.

6.1.3 Stationarity

Application of the SARIMAX model begins with ensuring data stationarity, which is a necessary condition for reliable time series forecasting. Using the Augmented Dickey-Fuller (ADF) test, each series in the data set is evaluated for stationarity. Also, the necessary order of differencing is determined to achieve this stationarity state. A summary table (Table 1) shows the variables and their respective differencing orders.

Variable	Differencing Order
TextBlob_SI	0
Vader_SI	1
Transformers_SI	0
Flair_SI	0
Afinn_SI	0
CCI	1
Average_wage	2
Pension	2
Inequality	2
Inflation	1
Unemployment_without_seasonality	1
Unemployment_rate	2

Table 1: Differencing Order Required for Stationarity

6.1.4 Selection of exogenous variables

The next step is to select exogenous variables that can improve the predictive ability of the model. This choice is based on Johansen’s cointegration test to identify variables that have a long-run equilibrium relationship with the CCI. Johansen’s test is typically used for time series that are I(1) (stationarity is obtained only after the first differencing) [22, 8]. Therefore, variables such as Vader_SI_diff_1, Inflation_diff_1, Unemployment_without_seasonality_diff_1, and 'CCI_diff_1' were selected for testing their cointegrated relationships (see Fig. 5). The test results suggest evidence of at least one and possibly more than one, but not more than two, cointegrating relationships among these four variables when considered together.

Lagged versions of selected variables are included in the model to capture temporal dependencies and lagged effects. Spearman’s rank correlation test, a nonparametric method that assesses the strength and direction of a monotonic relationship between two variables, is used to systematically determine optimal lagged variables.

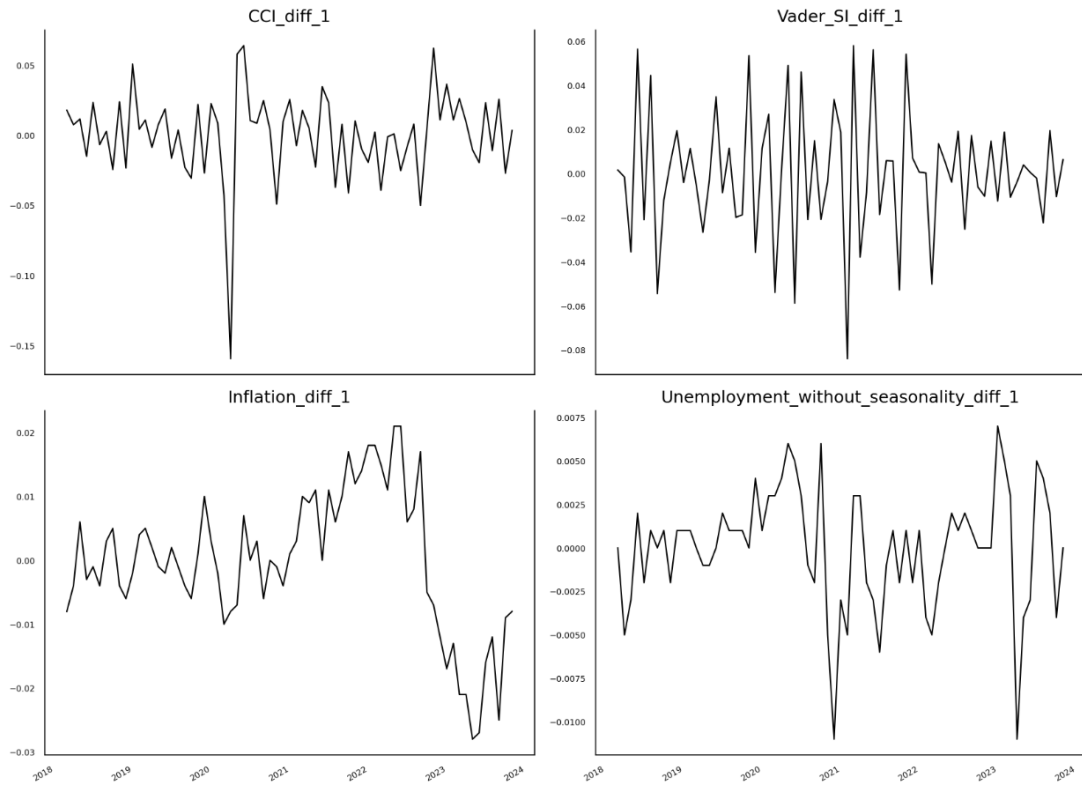


Figure 5: Differenced CCI and exogenous variables that have a long-run equilibrium relationship

The result of the correlation analysis is presented in a series of lines (Fig. 6), where each line shows the correlation of CCI_diff_1 with the exogenous variable at different lags. A lag of zero represents an unlagged variable, while positive lag values of 1, 2, 3, 6 and 12 months reflect the historical influence of exogenous variables on the CCI_diff_1.

It was decided to select Unemployment_without_seasonality_diff_1 with a six-month lag (0.22) and Inflation_diff_1 with its three-month lagged value (-0.29) which showed a moderate correlation with CCI_diff_1.

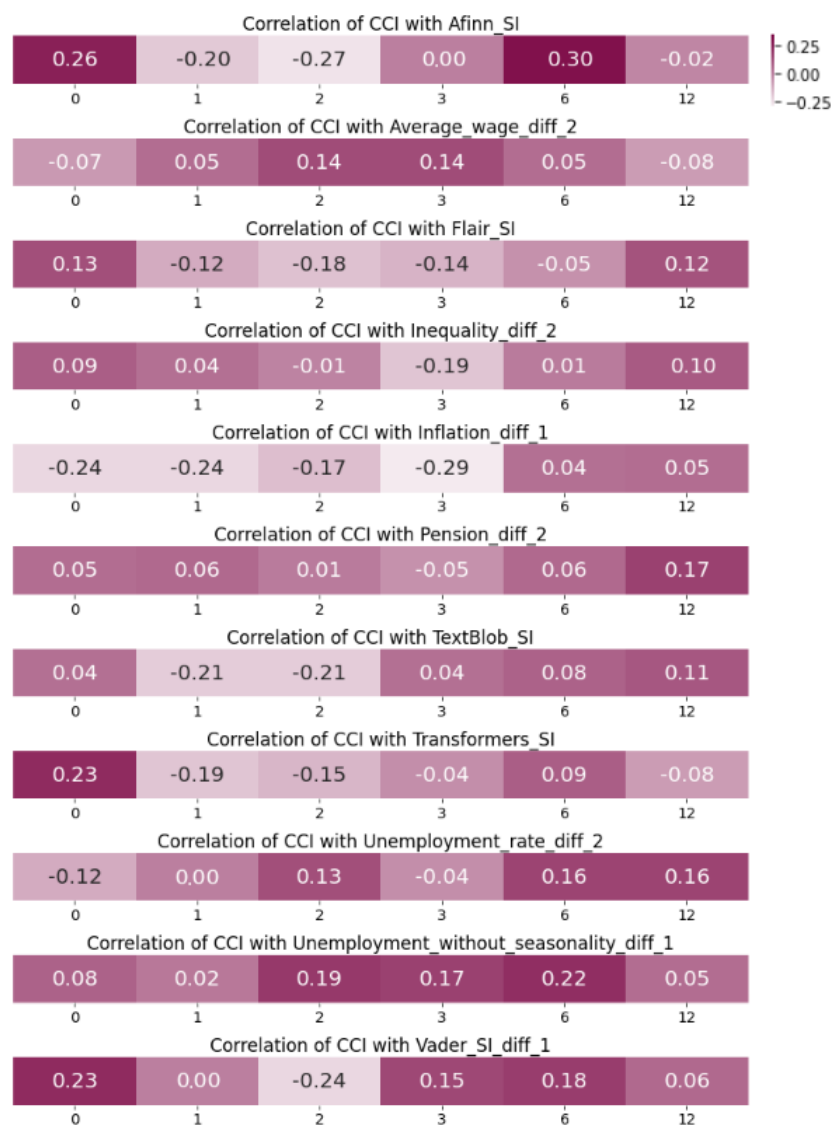


Figure 6: Correlation of CCI_diff_1 with exogenous variables and their lags

6.1.5 Optimal parameters selection using grid search

A grid search is performed to find the optimal parameters of the SARIMAX model. This involves iterating over a range of values for the non-seasonal and seasonal components of the SARIMAX(p,d,q)(P,D,Q)s model.

The order of integration d is set to 1 (series is stationary after being differenced once). The results of the ADF test on the seasonal component of CCI data show a very high negative ADF statistic and a p-value of 0.0. It means that the non-seasonal differencing (determined earlier as $d = 1$) is sufficient to make time series stationary, and no additional seasonal differencing is required for the SARIMAX model. Therefore $D = 0$. As data is recorded monthly, s is set to 12. The range up to 6 for 'p', 'q', 'P', and 'Q' is selected to capture the effects of recent trends and seasonal patterns up to a six-month horizon.

The optimal parameter combination for the SARIMAX model was selected based on MAE and AIC.

The MAE measures the model’s predictive accuracy, while the AIC provides insight into the model’s complexity, helping to avoid overfitting.

The combinations with parameters $p = 3$, $d = 1$, $q = 5$, $P = 0$, $D = 0$, and $Q = 3$ have an AIC of -222.366 and an MAE of 0.00812. This combination offers a good balance between a reasonably good AIC (while the lowest AIC was -232.469) and the lowest MAE. In this case, the emphasis is on MAE, because the model is used on an evolving dataset with new data points added monthly. Also, MAE is a comparative metric as it is used across other models in the study.

Therefore, the final model is SARIMAX(3,1,5)(0,0,3)12.

6.1.6 Residuals analysis

It is important to perform a residual analysis. Also, it is necessary to examine the Q-Q plot to ensure that it closely approximates a straight line, then apply the Ljung-Box test to see if the residuals are uncorrelated [34].

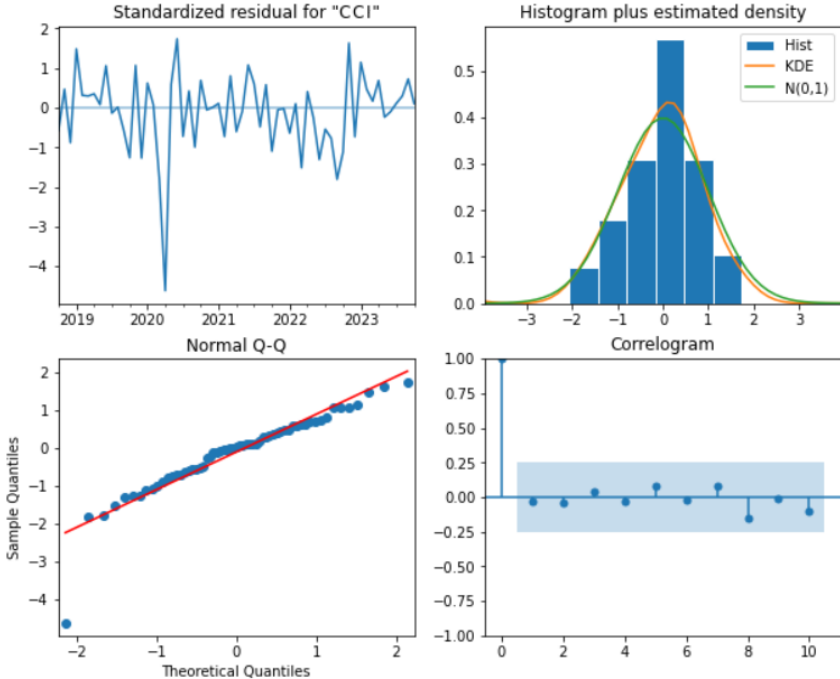


Figure 7: Diagnostic plots for testing normality and correlation of residuals

A diagnostic plots for testing normality and correlation of residuals is shown in Fig. 7. These plots are obtained after the final model fitting step. The standardized residuals plot shows no obvious systematic patterns or structural biases. The residuals oscillate around the zero line, indicating a well-fitted model that consistently captures the underlying trend of the time series without apparent bias. A close fit between the KDE (Kernel Density Estimate) and the normal distribution suggests that the residuals are approximately normally distributed. The alignment of data points with the theoretical line in a Normal Q-Q plot suggests that the residuals are normally distributed. The correlogram suggests that significant autocorrelation is present only at lag zero, indicating that the residuals are behaving as white noise.

The Ljung-Box test for the first 10 lags indicates all p-values are above the critical threshold of 0.05. It supports the null hypothesis that the residuals are independently distributed. Also, it confirms that the residuals are uncorrelated and the model’s forecasts are reliable.

6.1.7 Forecasting and results

Data for the last 7 values (approximately 10% of data) is forecasted using a rolling forecast approach, predicting one-time step at a time to minimize the accumulation of forecast errors (Fig. 8). This iterative process ensures that each forecast incorporates the most recent actual data, enhancing the model’s accuracy and responsiveness to new information.

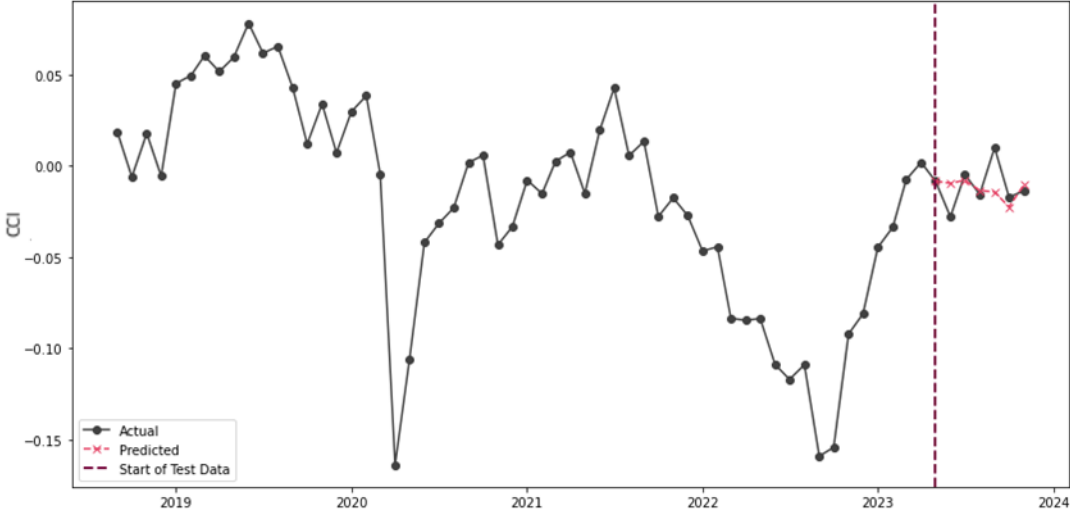


Figure 8: Actual vs predicted CCI SARIMAX(1,1,3)(1,0,3)12

6.2 VECM

6.2.1 VECM model fundamentals

The Vector Error Correction Model (VECM) is a model used for analyzing cointegrated time series [17]. It is particularly useful when dealing with non-stationary data that are integrated in the same order and cointegrated. Cointegration is essential for the VECM because it reflects the presence of long-run equilibrium relationships among variables. As noted by Engle et al. (1987) [15], cointegration indicates that certain sets of variables cannot significantly deviate from each other over time.

6.2.2 Model implementation

In developing the VECM for CCI forecasting, we used the main steps identified in the development of the SARIMAX model, taking into account the analogous data handling and variable selection requirements of both models. The main condition of both models is the stationarity of time series data. The Augmented Dickey-Fuller tests applied in the SARIMAX analysis to determine stationarity and to identify appropriate differences are also appropriate for the VECM model. The cointegration condition required for VECM is also verified during the analysis of the SARIMAX model.

The inclusion of exogenous variables in the SARIMAX model was based on Johansen’s cointegration test and Spearman’s rank correlation test. Based on these requirements, which are also relevant to the VECM model, we selected the same exogenous variables for the latter model.

Thus, we adopted a rolling forecasting approach, which is analogous to the method used in the SARIMAX model. We utilized the exogenous variables 'Vader_SI_diff_1', 'Inflation_diff_1', and 'Unemployment_without_seasonality_diff_1' in our model. Unlike the SARIMAX model, where we specified lagged versions of certain variables based on their correlation, the VECM inherently includes lagged versions of all exogenous variables. Therefore, to align with the nature of the VECM model, we chose a maximum lag of six (as for the SARIMAX model 'Unemployment_without_seasonality_diff_1' was lagged 6 months). However, the initial implementation of the model with these parameters resulted in a relatively high MAE of 0.14998, significantly larger than the MAE of 0.00812 achieved by the SARIMAX model.

6.2.3 Model optimization and results

To improve the model’s performance, we explored a systematic approach of hyperparameter tuning. This involved iterating over different combinations of exogenous variables and lag orders to identify the configuration that yielded the best forecast accuracy based on the MAE. We found that using only 'Unemployment_without_seasonality_diff_1' as the exogenous variable with a maximum lag of 4 resulted in a more acceptable MAE of 0.01597. Although this MAE was not as low as that of the SARIMAX model, it represented a significant improvement over our initial VECM implementation.

Moreover, we confirmed the presence of cointegrating relationships between 'CCI_diff_1' and 'Unemployment_without_seasonality_diff_1', which are essential for the application of VECM. The Johansen cointegration test results indicated at least two cointegrating relationships among these variables. This finding underscores the long-term equilibrium relationship between these variables.

The VECM model’s performance in forecasting the CCI is visually represented in the following figure (Fig. 9). It shows a comparison between the actual CCI values and those predicted by the VECM rolling forecasting method.

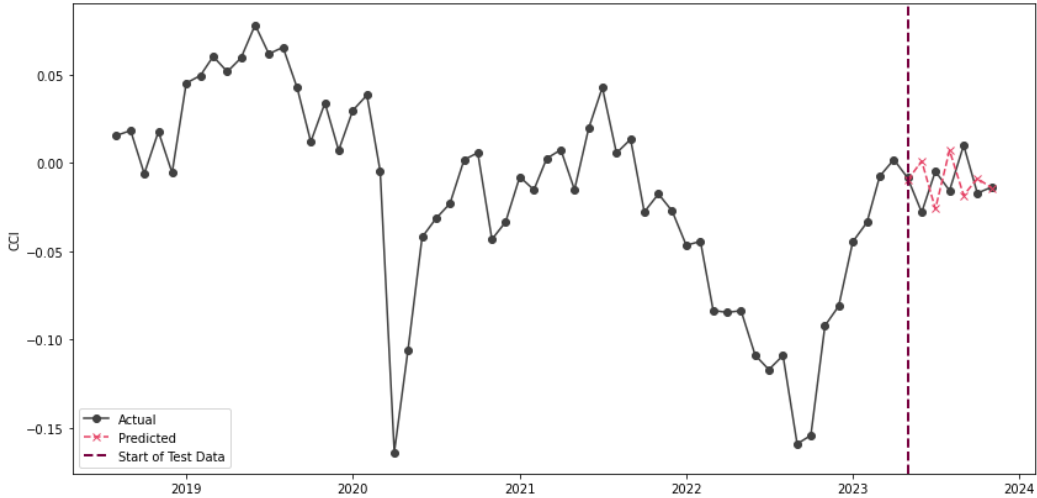


Figure 9: Comparison of actual vs. predicted CCI values using VECM rolling forecasting.

6.3 Random Forest

6.3.1 Random Forest model fundamentals

Random Forest (RF), known for its ability to handle complex datasets, is a powerful ensemble learning technique in machine learning. Breiman (2001) [7] explained that the Random Forest algorithm combines the outputs of numerous decision trees to get a single and more precise forecast. This helps avoid overfitting. Tyralis and Papacharalampous (2017) [38] explored RF for time series forecasting. They emphasised the importance of variable selection. Their study also examined how the number of lagged predictor variables affects the model’s accuracy. It used measures such as MAE, MSE and MAPE to compare the performance of different forecasting methods. Cutler et al. (2012) [13] mentioned one of the main advantages of RF - evaluate feature importance. It can automatically identify the most influential exogenous variables called features in a dataset.

6.3.2 Model implementation

To predict the CCI for the last 7 months, we used the rolling forecast method. At first, the model used all variables from the dataset, including their lags from 1 to 12 months to catch possible annual influence.

Additionally, hyperparameter tuning was conducted as part of the model training process to refine the model’s performance. This involved a systematic search for the best model parameters – such as the number of trees and their maximum depth – through a randomized search, which is a strategic approach to optimize model accuracy and minimize prediction errors.

After the initial model run (see Fig. 10a for the actual vs. predicted comparison), we get a quite low MAE = 0.01112 and the feature importance at each forecasting step (see Fig. 13). This helped identify the most influential features and their respective lags. Based on this analysis, the model was refined by selecting the top features and their specific lags for a more optimized forecasting model. The optimization process then is described in the following section.

6.3.3 Model optimization

We gained insights from the initial feature importance analysis in the model optimisation phase. The CCI lagged one month consistently held the most significant weight across multiple forecasting steps. This finding underscores the pivotal role that the CCI from the previous month plays in predicting the value for the ensuing month. Additionally, Inflation and Pension both lagged for 3 months and emerged as recurrently significant.

Further analysis revealed that incorporating sentiment indicators like TextBlob_SI in both unlagged and lagged forms (specifically, current, one-month lagged, and seven-month lagged) could enhance predictive performance. Applying this refined set of features resulted in a lower MAE, down to 0.01049 from the initial 0.01158 (see Fig. 10b). Feature importance after this step is shown in Fig. 14.

Taking it a step further, we incorporated rolling averages. This approach was predicated on the hypothesis that smoothing out short-term fluctuations could reveal underlying trends more pertinent

to the forecasting task. Adjusting the rolling windows to 4 months for Inflation and CCI, and 5 months for Pension, brings down the MAE to 0.00993 (see Fig. 10c, Fig. 15).

6.3.4 Model evaluation and results

The final CCI forecasting model for the past 7 months included a combination of non-lagged, lagged and rolling average features to optimize forecasting accuracy. The model features were CCI lagged by 1 month, TextBlob_SI in its current form, as well as lagged by 1 and 7 months, Inflation, and Pension lagged by 3 months. The rolling averages applied were 4 months for both Inflation and CCI, and 5 months for Pension.

The evaluation of the model’s performance and the illustrative results from each run can be observed in Fig. 10, where the graphical representations provide a comparison between the predicted and actual values of the CCI.

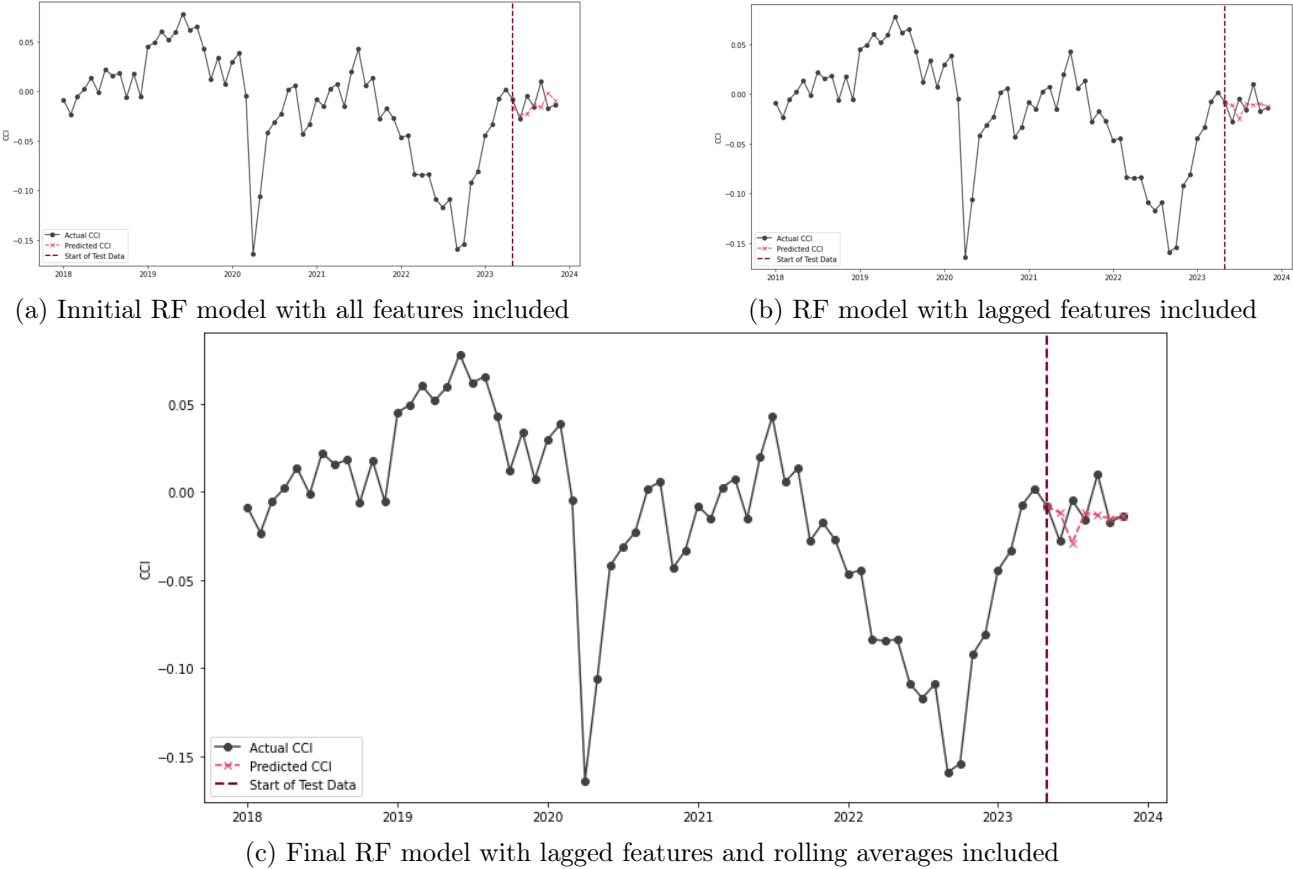


Figure 10: RF models

6.4 XGBoost

6.4.1 XGBoost model fundamentals

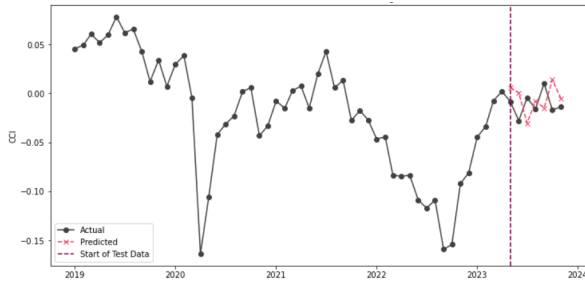
XGBoost (eXtreme Gradient Boosting) is an advanced implementation of gradient boosting algorithms. It has become a very popular tool for machine learning competitions. It can reveal complex non-linear patterns and common features of time series data. By analyzing time series, XGBoost builds an ensemble of regression trees using features such as early stopping to optimize training time and performance. A key feature of XGBoost is its ability to evaluate feature importance, which indicates how much each feature contributes to the prediction. This insight is particularly useful for identifying and eliminating non-contributing features from the model [9, 40].

6.4.2 Model implementation, optimization and results

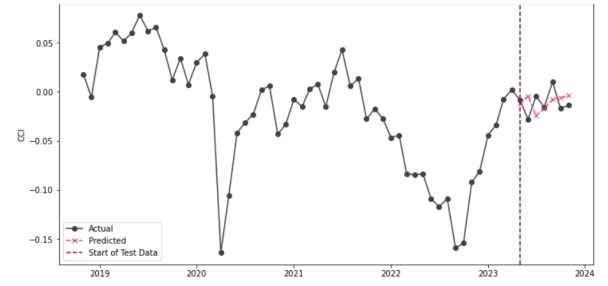
We employed a methodology akin to our approach with the Random Forest model. Initially, to forecast the last 7 months of CCI, we fed the model with all features from our dataset, incorporating their lagged values up to 12 months. Therefore, the MAE was relatively high at 0.01999 (see Fig. 11a for the actual vs. predicted comparison).

The feature importance graph (Fig. 16) allowed us to identify and focus on the most impactful features. Notably, we observed that including the CCI lagged by one month, the current value of 'Average_wage', 'Transformers_SI' shifted by two months, and 'Unemployment_without_seasonality' with a lag of 10 months enhanced the model's performance. To further investigate the influence of this latter variable, we also experimented with a shorter delay of just one month and with the current value. Our experimentation revealed that applying a 1-month lag and incorporating the current value (0-month lag) of this variable significantly improved the model's precision. Consequently, we chose to include only these short-term lags – 0 and 1 month – in our final model, which resulted in a lower MAE of 0.01396, as demonstrated in Fig. 11b, Fig. 17. This decision was based on the clear evidence that shorter lags provided a better forecast than the initially considered 10-month lag.

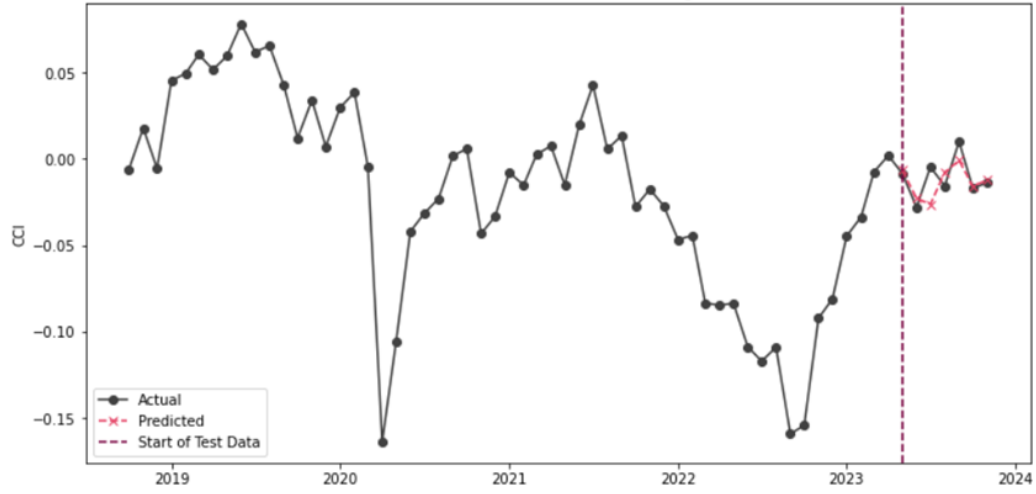
Also, we introduced rolling averages for key variables, utilizing a Rolling Window Grid Search for optimization. This approach involved testing various window sizes to find the best fit for each feature. For the CCI, a rolling window of 6 periods proved most effective. The average wage data benefited from a larger 9-period window, capturing longer-term trends. Sentiment analysis from the Transformers model and seasonally adjusted unemployment figures both showed optimal results with a 4-period window, allowing the model to quickly adapt to recent changes. Incorporating these specific rolling windows significantly improved our model's accuracy, reducing the MAE to 0.00713 (see Fig. 11c, Fig. 18).



(a) Initial model with all features included



(b) XGBoost model with lagged features included



(c) Final XGBoost model with lagged features and rolling averages included

Figure 11: XGBoost models

The improvements of our model after each optimization step are illustrated in Fig. 11, which shows a comparison of actual vs predicted values and demonstrates the increased accuracy achieved by systematically tuning efforts.

7 Conclusions

7.1 Summary of findings

In this work, we aimed to predict CCI. We explored 4 different models incorporating data from social media and administrative sources. The selection and integration of exogenous variables, taking into account their temporal influence through lags and rolling averages, helped to increase the models' accuracy. Table 2 shows the specific external variables used in each model.

The final models for SARIMAX, VECM, Random Forest, and XGBoost are presented in Fig. 12. This figure provides a visual comparison of the actual CCI values against the predictions generated by each model.

Model Name	Exogenous variables (features)
SARIMAX	Vader_SI_diff_1, Inflation_diff_1, Unemployment_without_seasonality_diff_1 (lag 6 months), Inflation_diff_1 (lag 3 months)
VECM	Unemployment_without_seasonality_diff_1 (lags up to 4 months) ¹
Random Forest	CCI (lag 1 month), TextBlob_SI (current, lag 1 and 7 months), Inflation (lag 3 months), Pension (lag 3 months), Rolling averages for Inflation (4 periods), CCI (4 periods), and Pension (5 periods)
XGBoost	CCI (lag 1 month), Average_wage (current), Transformers_SI (lag 2 months), Unemployment_without_seasonality (current and lag 1 month), Rolling averages for CCI (6 periods), Average_wage (9 periods), Unemployment (4 periods)

Table 2: Exogenous variables (features) in final forecasting models

The forecasting capabilities of each model were measured with Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and the Akaike Information Criterion (AIC) for SARIMAX. These metrics are presented in Table 3.

Based on the accuracy metrics (Table 3) and visual comparisons shown in the graphs (Fig. 12), XGBoost demonstrated better performance compared to the other models. The SARIMAX model followed closely, with a strong predictive ability, although slightly less accurate than XGBoost. The Random Forest and VECM models, while still valuable for their predictive insights, did not match the level of accuracy achieved by XGBoost and SARIMAX.

¹In the VECM model, $max_lags = 4$ applies to all variables in the system, meaning that up to 4 lagged values of each variable, including the variable CCI_diff_1, are included in the model.

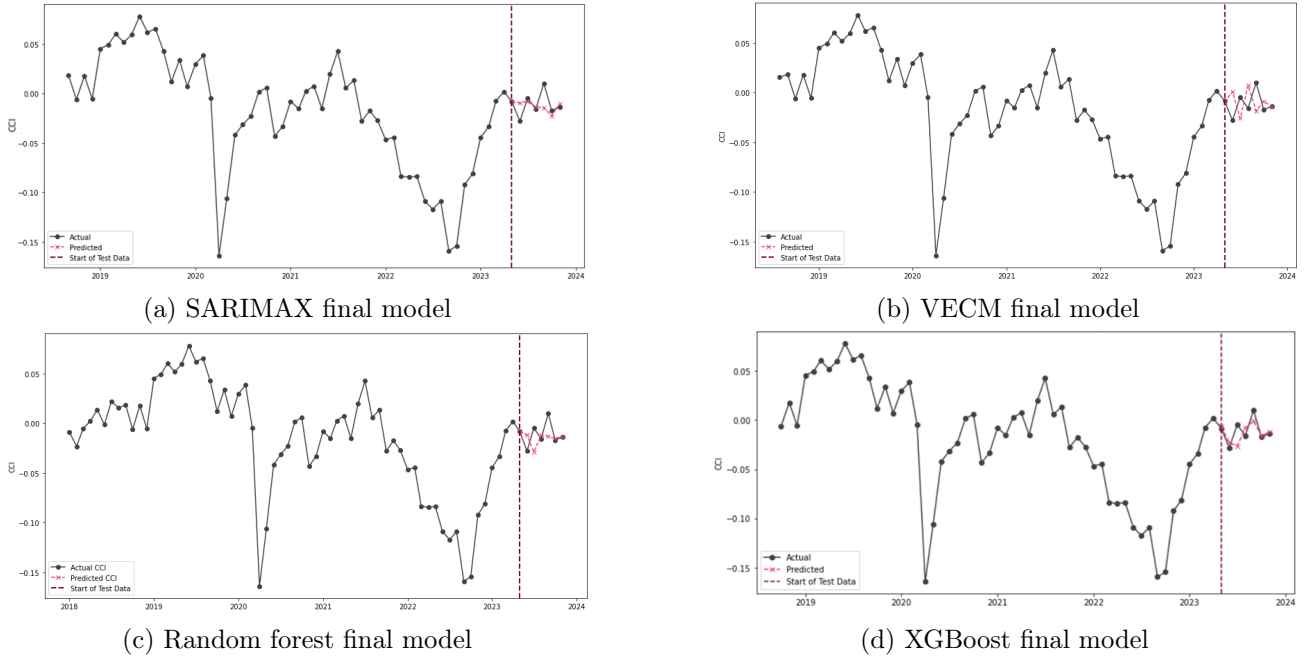


Figure 12: Final models for CCI forecasting

Model Name	MAE	MSE	RMSE
SARIMAX	0.00812	0.000066	0.00812
VECM	0.01597	0.000255	0.01597
Random Forest	0.00993	0.000098	0.00993
XGBoost	0.00713	0.000051	0.00713

Table 3: Forecasting errors of each model

7.2 Future work

The study offers a solid foundation for predicting of the Consumer Confidence Index (CCI) using social media and administrative data. However, it also opens up several prospects for future research. One such avenue is the exploration of advanced methodologies not yet explored in this thesis, such as Long Short Term Memory (LSTM) networks. These differ from the models currently employed. The thesis combines the econometric precision of SARIMAX and VECM models with the predictive capabilities of machine learning algorithms like XGBoost and Random Forest, which do not utilize the sequential memory characteristic of recurrent neural network (RNN) models like LSTM. Additionally, future research could benefit from incorporating and analyzing data from a variety of sources, including other economic indicators and Google Trends data to enhance the comprehensiveness and accuracy of the CCI prediction models.

References

- [1] R. Aishwarya, C. Ashwatha, A. Deepthi, and J. B. Raja. A novel adaptable approach for sentiment analysis. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, 5(2), 2019.
- [2] H. Akaike. Factor analysis and AIC. *Psychometrika*, 52:317–332, 1987.
- [3] Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. Flair: An easy-to-use framework for state-of-the-art nlp. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics (demonstrations)*, pages 54–59, 2019.
- [4] M. Ashraf, A. A. Raza, and M. Ishaq. A novel approach of social media analytics for predicting national consumer confidence index. *Bulletin of Business and Economics (BBE)*, 11(2):220–234, 2022.
- [5] E. Baldacci, B. Braaksma, A. G'alvez, K. Giannakouris, B. G. Olmos, P. Rivi'ere, M. Scannapieco, A. Vermeer, and V. Vertanen. Innovation during the covid-19 crisis: Why it was more critical for official statistics than ever. *Statistical Journal of the IAOS*, (Preprint):1–14, 2022.
- [6] J. van den Brakel, E. Söhler, P. Daas, and B. Buelens. Social media as a data source for official statistics; the dutch consumer confidence index. *Survey Methodology*, 43(2):183–210, 2017.
- [7] L. Breiman. Random Forests. *Machine learning*, 45:5–32, 2001.
- [8] E. Catanese, M. Scannapieco, M. Bruno, and L. Valentino. Natural language processing in official statistics: The social mood on economy index experience. *Statistical Journal of the IAOS*, (Preprint):1–9, 2022.
- [9] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [10] L. Corti, M. Zanetti, G. Tricella, and M. Bonati. Social media analysis of Twitter tweets related to asd in 2019–2020, with particular attention to covid-19: topic modelling and sentiment analysis. *Journal of big Data*, 9(1):113, 2022.
- [11] R. Curtin. Consumer sentiment surveys: worldwide review and assessment. *Journal of business cycle measurement and analysis*, 2007(1):7–42, 2007.
- [12] R. Curtin, S. Presser, and E. Singer. The effects of response rate changes on the index of consumer sentiment. *Public opinion quarterly*, 64(4):413–428, 2000.
- [13] D. R. Cutler, T. C. Edwards Jr, K. H. Beard, A. Cutler, K. T. Hess, J. Gibson, and J. J. Lawler. Random forests for classification in ecology. *Ecology*, 88(11):2783–2792, 2007.

- [14] S. K. Demirel and S. Artan. The causality relationships between economic confidence and fundamental macroeconomic indicators: Empirical evidence from selected european union countries. *International Journal of Economics and Financial Issues*, 7(5):417, 2017.
- [15] R. F. Engle and C. WJ. Granger. Co-integration and error correction: representation, estimation, and testing. *Econometrica: journal of the Econometric Society*, pages 251–276, 1987.
- [16] Eurostat, the statistical office of the European Union. Business and consumer surveys (source: Dg ecfm) (ei_bcs). https://ec.europa.eu/eurostat/cache/metadata/en/ei_bcs_esms.htm, 2023. Metadata last posted 30/05/2023.
- [17] S. Hajifar, H. Sun, F. M. Megahed, L. A. Jones-Farmer, E. Rashedi, and L. A. Cavuoto. A forecasting framework for predicting perceived fatigue: Using time series methods to forecast ratings of perceived exertion with features from wearable sensors. *Applied Ergonomics*, 90:103262, 2021.
- [18] E. Hargittai and E. Litt. The tweet smell of celebrity success: Explaining variation in Twitter adoption among a diverse group of young adults. *New media & society*, 13(5):824–842, 2011.
- [19] Clayton Hutto and Eric 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.
- [20] R. J. Hyndman and A. B. Koehler. Another look at measures of forecast accuracy. *International journal of forecasting*, 22(4):679–688, 2006.
- [21] T. U. Islam and M. N. Mumtaz. Consumer confidence index and economic growth: An empirical analysis of EU countries. *EuroEconomica*, 35(2), 2016.
- [22] S. Johansen. *Likelihood-based inference in cointegrated vector autoregressive models*. OUP Oxford, 1995.
- [23] G. Katona and R. Likert. Relationship between consumer expenditures and savings: The contribution of survey research. *The Review of Economics and Statistics*, 28(4):197–199, 1946.
- [24] S. Kemp. Digital 2023: Lithuania. <https://datareportal.com/reports/digital-2023-lithuania>, 2023. Published on February 13, 2023.
- [25] Y. W. Lee, K. G. Tay, and Y. Y. Choy. Forecasting electricity consumption using time series model. *International Journal of Engineering & Technology*, 7(4.30):218–223, 2018.
- [26] State Data Agency (Statistics Lithuania). Consumer survey. <https://osp.stat.gov.lt>, 2023. Metadata last posted 2023-05-29.
- [27] B. Liu. *Sentiment analysis and opinion mining*. Springer Nature, 2022.

- [28] J. Logemann. Measuring and managing expectations: Consumer confidence as an economic indicator, 1920s–1970s. *Futures Past. Economic Forecasting in the 20th and 21st Century*, 2020.
- [29] S. Loria et al. textblob documentation. *Release 0.16*, page 73, 2020.
- [30] P. Manigandan, M. S. Alam, M. Alharthi, U. Khan, K. Alagirisamy, D. Pachiyappan, and A. Rehman. Forecasting natural gas production and consumption in united states-evidence from sarima and sarimax models. *Energies*, 14(19):6021, 2021.
- [31] E. Mueller. Ten years of consumer attitude surveys: Their forecasting record. *Journal of the American Statistical Association*, 58(304):899–917, 1963.
- [32] F. Å. Nielsen. A new anew: Evaluation of a word list for sentiment analysis in microblogs. *arXiv preprint arXiv:1103.2903*, 2011.
- [33] L. Nowzohour and L. Stracca. More than a feeling: Confidence, uncertainty, and macroeconomic fluctuations. *Journal of Economic Surveys*, 34(4):691–726, 2020.
- [34] M. Peixeiro. *Time series forecasting in Python*. Simon and Schuster, 2022.
- [35] S. Stieglitz, M. Mirbabaie, B. Ross, and C. Neuberger. Social media analytics—challenges in topic discovery, data collection, and data preparation. *International journal of information management*, 39:156–168, 2018.
- [36] The Conference Board. Consumer confidence survey technical note – may 2021. <https://www.conference-board.org/>, May 2021.
- [37] A. Tumasjan, T. Sprenger, P. Sandner, and I. Welpe. Predicting elections with Twitter: What 140 characters reveal about political sentiment. In *Proceedings of the international AAAI conference on web and social media*, volume 4, pages 178–185, 2010.
- [38] H. Tyrallis and G. Papacharalampous. Variable selection in time series forecasting using Random Forests. *Algorithms*, 10(4):114, 2017.
- [39] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.
- [40] N. Woloszko. Adaptive trees: a new approach to economic forecasting. 2020.

Appendices

Išlaidos	Išlaidas	Išlaidų	Išleidžiu	Išleidžiamė
Išleisti	Kaštai	Santaupos	Sutaupau	Santaupų
Santaupas	Taipome	Taupau	Taupyti	Taupymas
Sutaupome	Sutaupyti	Pirkti	Pirkinys	Pirkiniai
Pirkinių	Perka	Perku	Pirksite	Perkame
Nuperku	Nusiperku	Bedarbiai	Bedarbis	Bedarbių
Bedarbiams	Bedarbius	Nedarbas	Nedirbu	Nedirbti
Darbas	Darbai	Darbą	Darbu	Dirbti
Dirbu	Ekonomika	Ekonominė padėtis	Finansinė padėtis	Finansai
Finansus	Finansinis	Finansų	Pinigai	Pinigus
Pinigų	Pinigais	Kainuoti	Kainuoja	Kaina
Kainos	Kainose	Kainų	Kainų pokyčiai	Defliacija
Infliacija	Inflacijai	Infliuoti	Infliacijos	Atlyginimas
Atlyginima	Atlyginimai	Alga	Algos	Pajamos
Pajamomis	Pajamų	Palūkanų normos	Palūkanų	Palūkanos
Palūkanas	Biudžetas	Biudžetą	Biudžetui	Biudžetu
BVP	Bendras vidaus produktas	Rinka	Mokesčiai	Mokesčius
Mokestis	Sumokėti	Mokėjimas	Pensija	Pensijos
Pensijas	Skola	Skolintis	Skolinti	Skolinu
Skolinuosi	Pasiskolinti	Įsiskolinimas	Įsiskolinu	Įsiskolinti
Namų ūkis				

Table 4: List of economic keywords

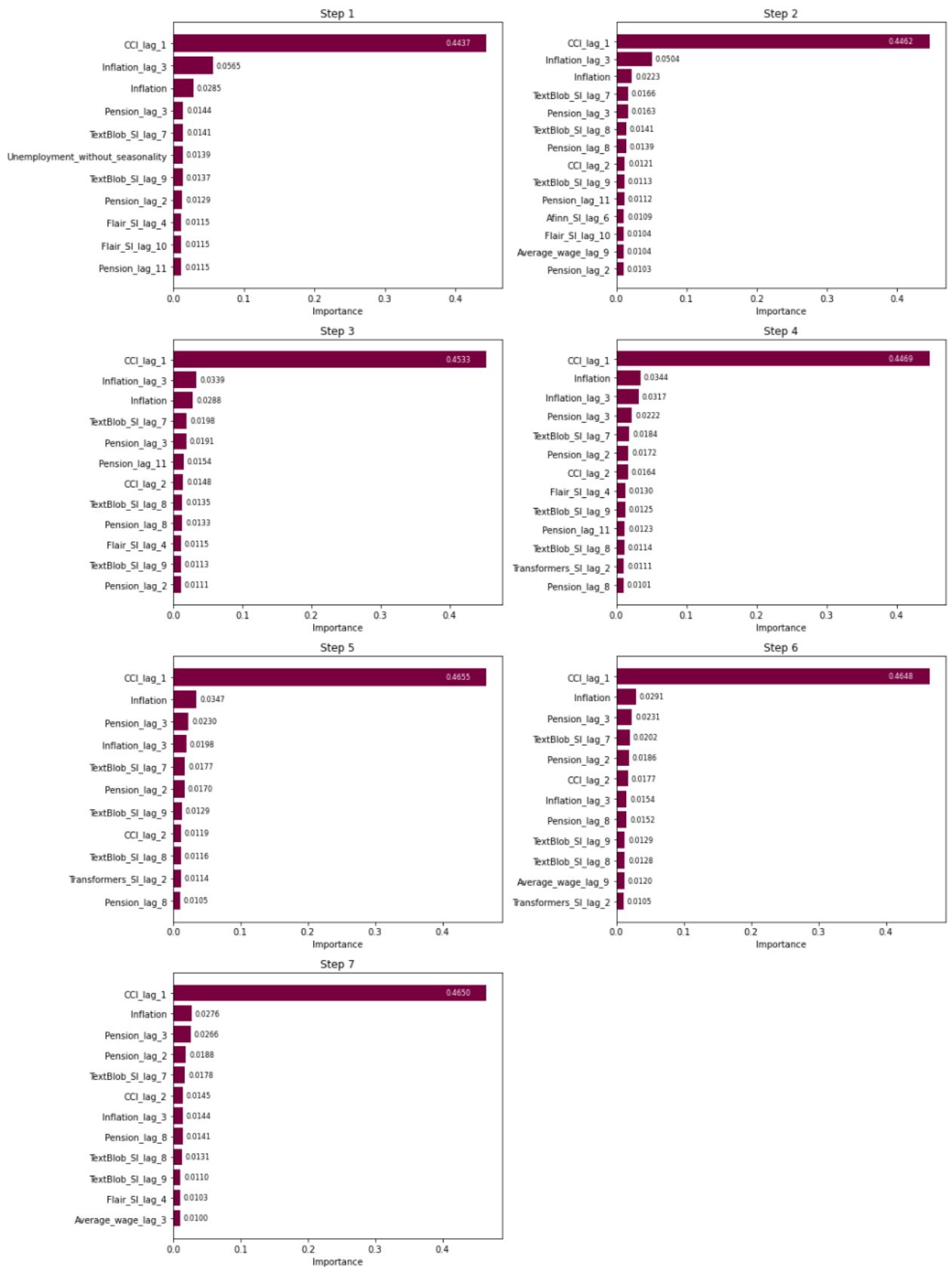


Figure 13: Feature importance for Random Forest model using all features

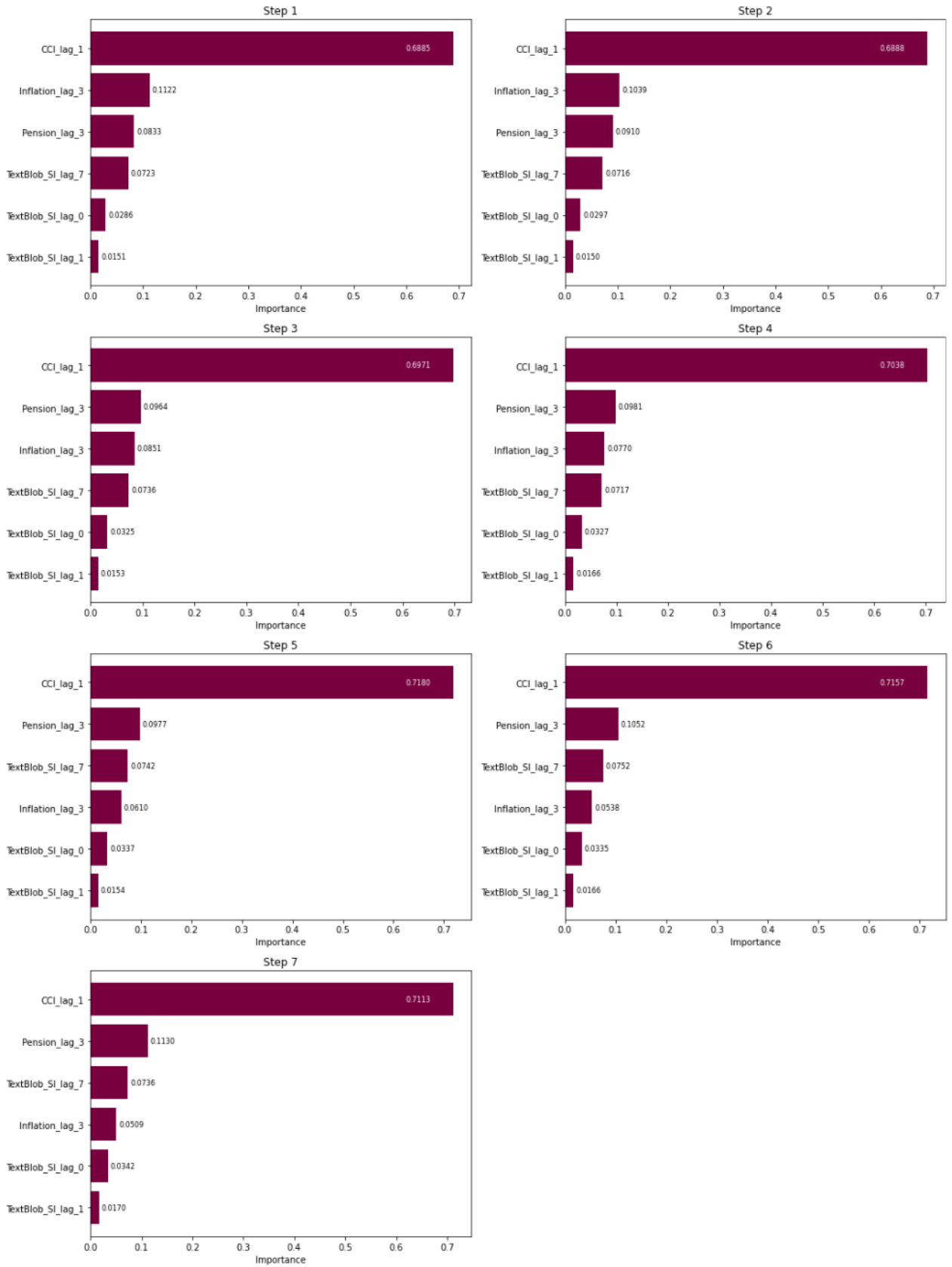


Figure 14: Feature importance for Random Forest model using selected features and their lags

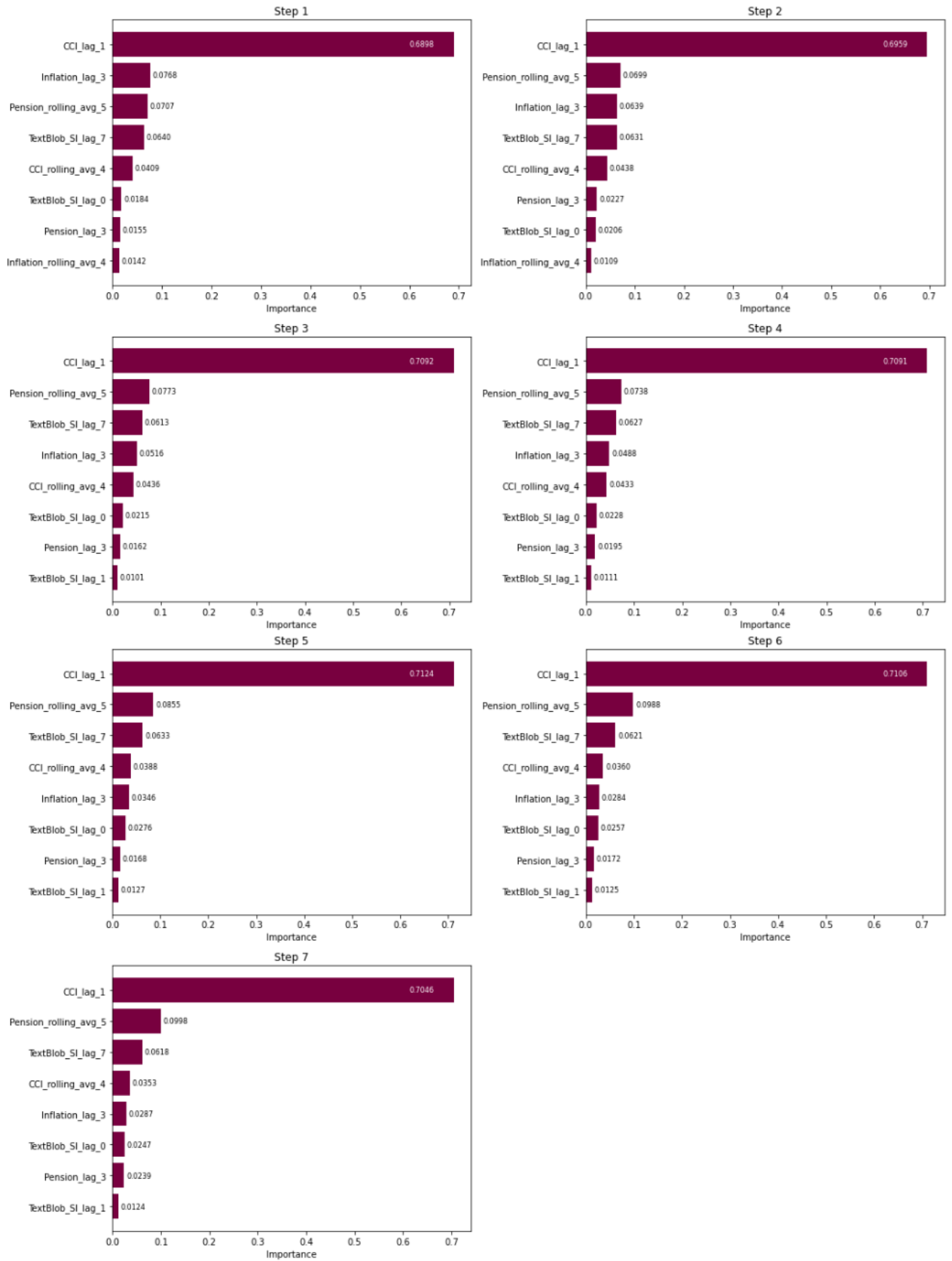


Figure 15: Feature importance for Random Forest model using selected features, their lags and rolling averages

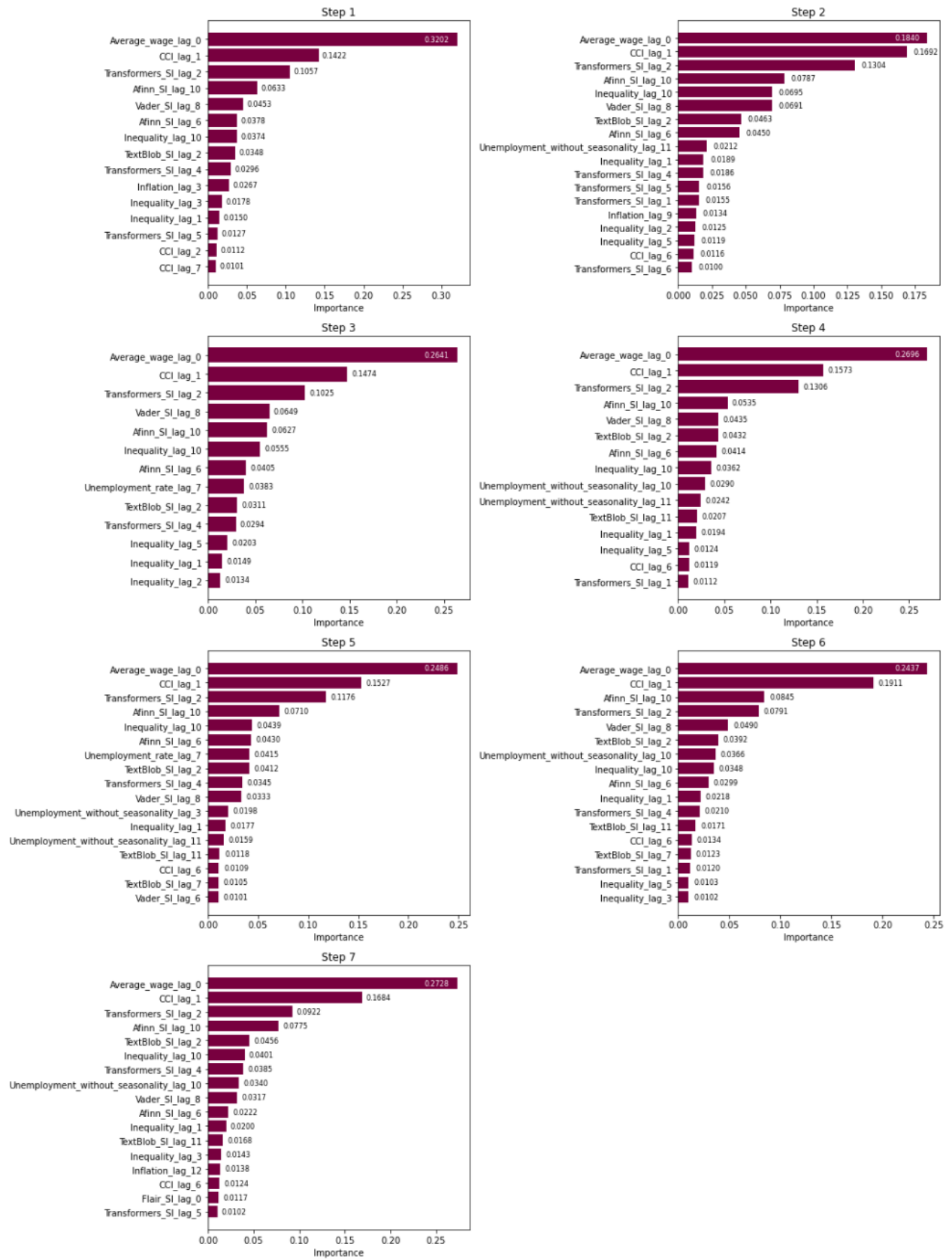


Figure 16: Feature importance for XGBoost model using all features

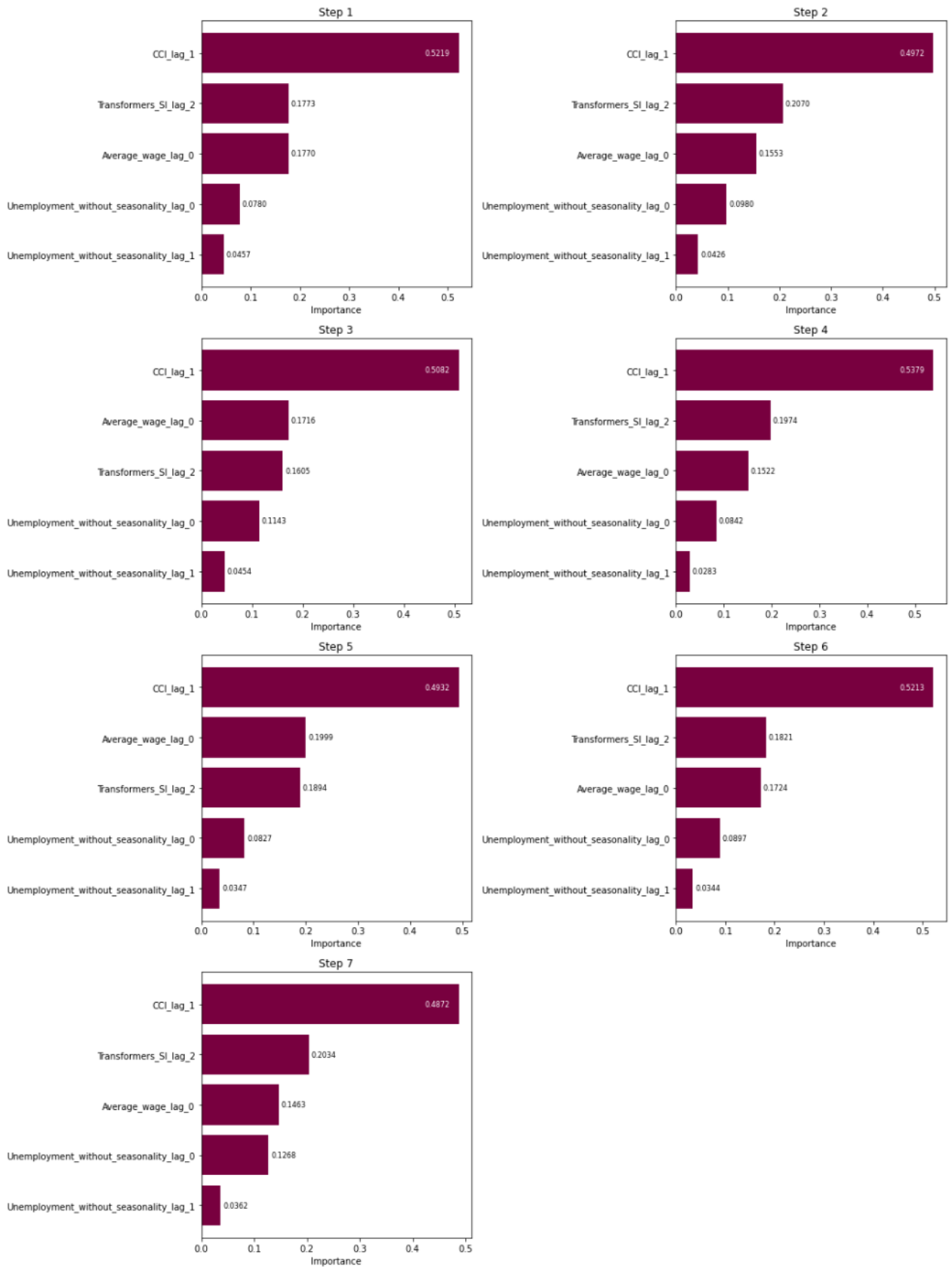


Figure 17: Feature importance for XGBoost model using selected features and their lags

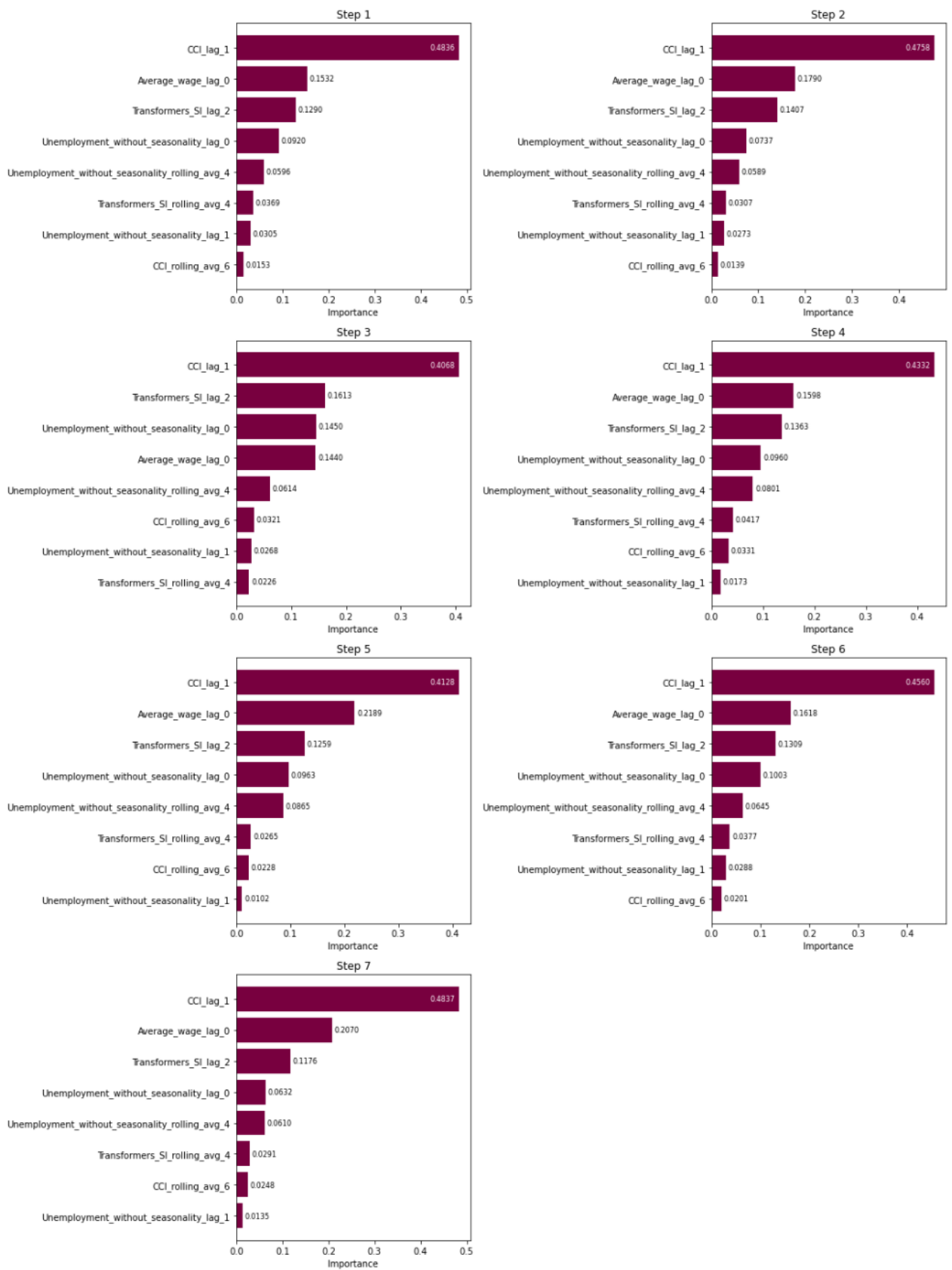


Figure 18: Feature importance for XGBoost model using selected features, their lags and rolling averages

Here we provide the main code parts used.

```
1 #— Libraries —————
2 # Data Handling
3 import pandas as pd
4 import numpy as np
5 import os
6 import os.path
7 import time
8 import datetime as dt
9 from datetime import datetime
10
11 # Text and Language Processing
12 from unidecode import unidecode
13 from langdetect import detect
14 import re
15
16 # Twitter API Interaction
17 import tweepy
18
19 # Sentiment Analysis Libraries
20 from textblob import TextBlob
21 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
22 from transformers import pipeline
23 import flair
24 import stanza
25 from afinn import Affin
26
27 # File Handling
28 import xlswriter
29 import openpyxl
30
31 # HTTP and SSL for API Interaction
32 import requests
33 from requests.adapters import HTTPAdapter
34 from urllib3.util.ssl_ import create_urllib3_context
35 import ssl
36
37 # Statistical and Forecasting Models
38 from statsmodels.tsa.seasonal import STL
39 from statsmodels.tsa.statespace.sarimax import SARIMAX
40 from statsmodels.tsa.vector_ar.vecm import VECM
41 from statsmodels.tsa.stattools import adfuller
42 import scipy.stats as stats
43
44 # Machine Learning Models and Evaluation
45 from sklearn.ensemble import RandomForestRegressor
46 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
47 from sklearn.metrics import mean_squared_error, mean_absolute_error,
    mean_absolute_percentage_error
```



```

48 import xgboost as xgb
49
50 # Visualization
51 import matplotlib.pyplot as plt
52 import seaborn as sns
53
54 # Miscellaneous
55 from itertools import product
56 from math import sqrt
57 from pprint import pprint
58
59 #-----Weekly scraping-----
60 client = tweepy.Client(bearer_token=bearer_token)
61 # tweet_fields
62 tweet_fields = [ 'attachments',
63                 'author_id',
64                 'context_annotations',
65                 'conversation_id',
66                 'created_at',
67                 'edit_controls',
68                 'edit_history_tweet_ids',
69                 'entities',
70                 'geo',
71                 'id',
72                 'in_reply_to_user_id',
73                 'lang',
74                 'possibly_sensitive',
75                 'public_metrics',
76                 'referenced_tweets',
77                 'reply_settings',
78                 'source',
79                 'text',
80                 'withheld'
81                 ]
82
83 user_fields = [ 'created_at',
84               'description',
85               'entities',
86               'id',
87               'location',
88               'pinned_tweet_id',
89               'profile_image_url',
90               'protected',
91               'public_metrics',
92               'url',
93               'verified',
94               'verified_type',
95               'withheld'
96               ]

```

```

97
98 place_fields = [ 'contained_within',
99                 'country',
100                'country_code',
101                'full_name',
102                'geo',
103                'id',
104                'name',
105                'place_type'
106                ]
107
108 expansions = [ 'author_id',
109               'geo.place_id',
110               'attachments.media_keys'
111               ]
112
113 def search_tweets(query, tweet_fields, user_fields, place_fields, expansions,
114                  start_time, max_results):
115     """
116     This function searches for recent tweets based on various parameters.
117
118     Parameters:
119     - query: The search query to find tweets.
120     - tweet_fields: Fields to include in the tweet objects.
121     - user_fields: Fields to include in the user objects.
122     - place_fields: Fields to include in the place objects.
123     - expansions: Expansions to include additional objects in the payload.
124     - start_time: The oldest UTC timestamp from which the tweets will be provided.
125     - max_results: The max number of search results returned by a single API call.
126
127     Returns:
128     A list of tweets that match the search criteria.
129     """
130     tweets = client.search_recent_tweets(query = query,
131                                         tweet_fields = tweet_fields,
132                                         user_fields = user_fields,
133                                         place_fields = place_fields,
134                                         expansions = expansions,
135                                         start_time = start_time,
136                                         max_results = max_results)
137     return tweets
138
139 def get_tweets_df(queries, tweet_fields, user_fields, place_fields, expansions,
140                  start_time, max_results, name_ex, file_handle=None):
141     """
142     Retrieves and saves tweet data for a given set of queries.
143     The function processes each query, retrieves relevant tweets, and saves the

```

```

144     tweet information in an Excel file. It logs each step and the total number of
145     tweets scraped at the end of the process.
146     """
147     tweet_data = []
148     df_tweets = []
149     places = {}
150     users = {}
151     i = 0
152     num_rows = 0
153     start_position = 0
154     total_tweets = 0
155
156     notebook_output = ""
157
158     for query in queries:
159         notebook_output += f"{i}\n"
160
161         started_time = datetime.now().strftime('%Y-%m-%d %H:%M:%S')
162         notebook_output += f"Processing query: '{query}' | Start Time: {started_time}\n"
163         "
164         tweets_query = search_tweets(query=query + str(" -is:reply -is:retweet lang:lt"
165         ),
166                                     tweet_fields=tweet_fields,
167                                     user_fields=user_fields,
168                                     place_fields=place_fields,
169                                     expansions=expansions,
170                                     start_time = start_time,
171                                     #end_time = end_time,
172                                     max_results=max_results)
173
174         if tweets_query.data is not None:
175             # Save tweets in a text file
176             path = txt_tweet_path + today_date + '_' + query + '.txt'
177             with open(path, 'w', encoding='utf-8') as file:
178                 file.write(str(tweets_query))
179
180             if 'places' in tweets_query.includes:
181                 places.update({p["id"]: p for p in tweets_query.includes['places']})
182
183             if 'users' in tweets_query.includes:
184                 users.update({u["id"]: u for u in tweets_query.includes['users']})
185
186             for tweet in tweets_query.data:
187                 tweet_info = {
188                     'Query': query,
189                     'Tweet_ID': tweet.id,
190                     'Author_ID': tweet.author_id,

```

```

190         'Text': tweet.text,
191         'Created_At': tweet.created_at,
192         'Language': tweet.lang,
193         'Public_Metrics': tweet.public_metrics,
194         'Attachments': tweet.attachments,
195         'Context_Annotations': tweet.context_annotations,
196         'Conversation_ID': tweet.conversation_id,
197         'Edit_Controls': tweet.edit_controls,
198         'Edit_History_Tweet_IDs': tweet.edit_history_tweet_ids,
199         'Entities': tweet.entities,
200         'In_Reply_To_User_ID': tweet.in_reply_to_user_id,
201         'Possibly_Sensitive': tweet.possibly_sensitive,
202         'Referenced_Tweets': tweet.referenced_tweets,
203         'Reply_Settings': tweet.reply_settings,
204         'Source': tweet.source,
205         'Withheld': tweet.withheld,
206         'Geo': tweet.geo
207     }
208
209     if tweet.geo is not None and tweet.geo['place_id'] in places:
210         place = places[tweet.geo['place_id']]
211         tweet_info['Geo_Name'] = place['full_name']
212         tweet_info['Country'] = place['country']
213         tweet_info['Short_Geo_Name'] = place['name']
214         tweet_info['Place_type'] = place['place_type']
215
216     else:
217         tweet_info['Geo_Name'] = ""
218         tweet_info['Country'] = ""
219         tweet_info['Short_Geo_Name'] = ""
220         tweet_info['Place_type'] = ""
221
222     if tweet.author_id is not None and tweet.author_id in users:
223         user = users[tweet.author_id]
224         tweet_info['Location'] = user['location']
225         #tweet_info['User_ID'] = user['id']
226
227     else:
228         tweet_info['Location'] = ""
229         #tweet_info['User ID'] = ""
230
231     tweet_info['Timestamp'] = started_time
232     tweet_data.append(tweet_info)
233
234 df_tweets = pd.DataFrame(tweet_data)
235
236 # Convert datetime column to a specific timezone (e.g., UTC)
237 df_tweets['Created_At'] = pd.to_datetime(df_tweets['Created_At'])\
238     .dt.tz_convert('UTC')

```

```

239
240     # Convert datetime column to string format without timezone information
241     df_tweets['Created_At'] = df_tweets['Created_At'].dt\
242         .strftime('%Y-%m-%d %H:%M:%S')
243
244     if i==0:
245         start_position = 0
246     else:
247         start_position = start_position + num_rows
248
249     # create the excel file if it not already exists
250     if not os.path.exists(excel_tweets_path+"Tweets_"+today_date+name_ex+'.xlsx
251 ):
252         workbook = xlswriter.Workbook(excel_tweets_path+"Tweets_"
253         +today_date+name_ex+'.xlsx')
254         workbook.close()
255
256     # write data in the excel file with name based on search_phrase_combination
257     with pd.ExcelWriter(excel_tweets_path+"Tweets_"+today_date+name_ex+'.xlsx',
258         engine="openpyxl",
259         if_sheet_exists='overlay',
260         mode='a') as writer:
261         df_tweets.to_excel(writer,
262             startrow=start_position,
263             startcol=0,
264             index=False,
265             header=False)
266
267     num_rows = len(df_tweets)
268     #print(f"Processing query: '{query}' | Tweets: {num_rows}")
269     notebook_output += f"Processing query: '{query}' | Tweets: {num_rows}\n"
270     total_tweets += num_rows
271
272     ended_time = datetime.now().strftime('%Y-%m-%d %H:%M:%S')
273     #print(f"Processing query: '{query}' | End Time: {ended_time}\n")
274     notebook_output += f"Processing query: '{query}' | End Time: {ended_time}\n
275 \n"
276
277     else:
278         #print("No tweet data available for the query. \n")
279         notebook_output += "No tweet data available for the query.\n\n"
280
281     #num_rows = 0
282     df_tweets = []
283     tweet_data = []
284     places = {}
285     users = {}
286     i = i + 1

```

```

286 # Add the total at the end
287 notebook_output += f"Total Tweets Scraped: {total_tweets}\n"
288
289 # Print the output to the notebook
290 print(notebook_output)
291
292 # If a file handle is provided, save the output to the text file
293 if file_handle:
294     file_handle.write(notebook_output)
295     file_handle.flush()
296
297 #-----Timelines scraping-----
298
299 class AdapterFix(HTTPAdapter):
300     def __init__(self):
301         self.ssl_context = ssl.create_default_context()
302         super().__init__()
303
304     def init_poolmanager(self, *args, **kwargs):
305         kwargs["ssl_context"] = self.ssl_context
306         super().init_poolmanager(*args, **kwargs)
307
308     def proxy_manager_for(self, *args, **kwargs):
309         kwargs["ssl_context"] = self.ssl_context
310         return super().proxy_manager_for(*args, **kwargs)
311
312 def Client(bearer_token):
313     adapter = AdapterFix()
314
315     client = tweepy.Client(bearer_token=bearer_token)
316     client.session.adapters['https://'] = adapter
317     # this needs to be done before calling any API methods.
318     return client
319
320 def get_today():
321     # Get today's date
322     today = dt.date.today()
323     # Format the date as YY_MM_DD
324     today_date = today.strftime('%y_%m_%d')
325     print(today_date)
326     return(today_date)
327
328 def get_users_tweets_df(ids, folder, tweet_fields, user_fields, place_fields,
329     expansions, exclude, start_time, end_time, max_results, name_ex, k, bearer_token,
330     file_handle=None):
331     today_date = get_today()
332     tweet_data = []
333     df_tweets = []

```

```

333 places = {}
334 users = {}
335 i = 0
336 num_rows = 0
337 start_position = 0
338 total_tweets = 0
339
340 notebook_output = ""
341 notebook_output += f"Beginning of the period | Time: {start_time}\n\n"
342
343 for id in ids:
344
345     #if i == 0:
346
347     notebook_output += f"{id}\n"
348
349     started_time = datetime.now().strftime('%Y-%m-%d %H:%M:%S')
350     notebook_output += f"Processing id: '{id}' | Start Time: {started_time}\n"
351
352     tweets_id = get_users_tweets(id=id,
353                                 tweet_fields=tweet_fields,
354                                 user_fields=user_fields,
355                                 place_fields=place_fields,
356                                 expansions=expansions,
357                                 exclude = exclude,
358                                 start_time = start_time,
359                                 end_time = end_time,
360                                 max_results=max_results,
361                                 bearer_token=bearer_token)
362
363     if tweets_id.data is not None:
364         # Save tweets in a text file
365         path = txt_tweet_path + today_date + '_' + id + '.txt'
366         with open(path, 'w', encoding='utf-8') as file:
367             file.write(str(tweets_id))
368
369         if 'places' in tweets_id.includes:
370             places.update({p["id"]: p for p in tweets_id.includes['places']})
371
372         if 'users' in tweets_id.includes:
373             users.update({u["id"]: u for u in tweets_id.includes['users']})
374
375     for tweet in tweets_id.data:
376         tweet_info = {
377             'ID': id,
378             'Tweet_ID': tweet.id,
379             'Author_ID': tweet.author_id,
380             'Text': tweet.text,
381             'Created_At': tweet.created_at,

```

```

382         'Language': tweet.lang,
383         'Public_Metrics': tweet.public_metrics,
384         'Attachments': tweet.attachments,
385         'Context_Annotations': tweet.context_annotations,
386         'Conversation_ID': tweet.conversation_id,
387         'Edit_Controls': tweet.edit_controls,
388         'Edit_History_Tweet_IDs': tweet.edit_history_tweet_ids,
389         'Entities': tweet.entities,
390         'In_Reply_To_User_ID': tweet.in_reply_to_user_id,
391         'Possibly_Sensitive': tweet.possibly_sensitive,
392         'Referenced_Tweets': tweet.referenced_tweets,
393         'Reply_Settings': tweet.reply_settings,
394         'Source': tweet.source,
395         'Withheld': tweet.withheld,
396         'Geo': tweet.geo
397     }
398
399     if tweet.geo is not None and tweet.geo['place_id'] in places:
400         place = places[tweet.geo['place_id']]
401         tweet_info['Geo_Name'] = place['full_name']
402         tweet_info['Country'] = place['country']
403         tweet_info['Short_Geo_Name'] = place['name']
404         tweet_info['Place_type'] = place['place_type']
405
406     else:
407         tweet_info['Geo_Name'] = ""
408         tweet_info['Country'] = ""
409         tweet_info['Short_Geo_Name'] = ""
410         tweet_info['Place_type'] = ""
411
412     if tweet.author_id is not None and tweet.author_id in users:
413         user = users[tweet.author_id]
414         tweet_info['Location'] = user['location']
415         #tweet_info['User_ID'] = user['id']
416
417     else:
418         tweet_info['Location'] = ""
419         #tweet_info['User ID'] = ""
420
421     tweet_info['Timestamp'] = started_time
422
423     tweet_data.append(tweet_info)
424
425 df_tweets = pd.DataFrame(tweet_data)
426
427 # Convert datetime column to a specific timezone (e.g., UTC)
428 df_tweets['Created_At'] = pd.to_datetime(df_tweets['Created_At'])\
429     .dt.tz_convert('UTC')
430

```



```

431 # Convert datetime column to string format without timezone information
432 df_tweets['Created_At'] = df_tweets['Created_At']\
433     .dt.strftime('%Y-%m-%d %H:%M:%S')
434
435 if i==0:
436     start_position = 0
437 else:
438     start_position = start_position + num_rows
439
440 # create the excel file if it not already exists
441 if not os.path.exists(excel_user_tweets_path
442     +folder
443     +"User_tweets_"
444     +today_date+"_"
445     +str(k)+name_ex
446     +'.xlsx'):
447     workbook = xlsxwriter.Workbook(excel_user_tweets_path
448     +folder
449     +"User_tweets_"
450     +today_date+"_"
451     +str(k)+name_ex+'.xlsx')
452     workbook.close()
453
454 # write data in the excel file with name based on search_phrase_combination
455 with pd.ExcelWriter(excel_user_tweets_path+folder
456     +"User_tweets_"
457     +today_date
458     +"_"
459     +str(k)
460     +name_ex
461     +'.xlsx',
462     engine="openpyxl",
463     if_sheet_exists='overlay',
464     mode='a') as writer:
465     df_tweets.to_excel(writer,
466         startrow=start_position,
467         startcol=0,
468         index=False,
469         header=False)
470
471 num_rows = len(df_tweets)
472 #print(f"Processing id: '{id}' | Tweets: {num_rows}")
473 notebook_output += f"Processing id: '{id}' | Tweets: {num_rows}\n"
474 total_tweets += num_rows
475
476 ended_time = datetime.now().strftime('%Y-%m-%d %H:%M:%S')
477 #print(f"Processing id: '{id}' | End Time: {ended_time}\n")
478 notebook_output += f"Processing id: '{id}' | End Time: {ended_time}\n\n"
479

```

```

480     else:
481         #print("No tweet data available for this user. \n")
482         notebook_output += "No tweet data available for this user.\n\n"
483
484         df_tweets = []
485         tweet_data = []
486         places = {}
487         users = {}
488         i = i + 1
489
490     #print(f"Enf of the period | Time: {end_time}")
491     notebook_output += f"End of the period | Time: {end_time}\n"
492     notebook_output += f"Total Tweets Scraped: {total_tweets}\n" # Add the total at
the end
493
494     # Print the output to the notebook
495     print(notebook_output)
496
497     # If a file handle is provided, save the output to the text file
498     if file_handle:
499         file_handle.write(notebook_output)
500         file_handle.flush()
501
502     return total_tweets
503
504 def call_function_with_breaks(ids, folder, start_time, end_time, bearer_token, kk):
505     today_date = get_today()
506     chunk_size = 5 # Number of IDs to process in each iteration
507     maxi = 100 #tweets for id
508     total_tweets_all_iterations = 0
509
510     # Create and open the text file to save the output
511     file_path = excel_user_tweets_path+folder
512                 +str(folder)[: -1]
513                 +"_LOG_"
514                 +today_date
515                 +"_"
516                 +str(kk)
517                 +".txt"
518     with open(file_path, "w") as file_handle:
519         # Iterate over the IDs in chunks of size 'chunk_size'
520         for i in range(0, len(ids), chunk_size):
521             # Get the current chunk of IDs
522             chunk = ids[i:i+chunk_size]
523
524             # Call the function with the current chunk of IDs
525             total_tweets_iteration = get_users_tweets_df(ids=chunk,
526                                                         folder=folder,
527                                                         tweet_fields=tweet_fields,

```

```

528         user_fields=user_fields ,
529         place_fields=place_fields ,
530         expansions=expansions ,
531         exclude=exclude ,
532         start_time=start_time ,
533         end_time=end_time ,
534         max_results=maxi ,
535         name_ex = "_" + str(int(i/4)) ,
536         k = kk ,
537         bearer_token = bearer_token ,
538         file_handle=file_handle)
539
540         total_tweets_all_iterations += total_tweets_iteration
541         print(i+5)
542         print(f"Total Tweets Scraped: {total_tweets_all_iterations}")
543         time.sleep(15.1 * 60) # Sleep for 15.1 minutes in seconds
544     file_handle.close() # Close the file after the with block ends
545
546     # Open the text file and write the final total_tweets_all_iterations
547     with open(file_path, "a") as file_handle:
548         file_handle.write(f"Final Total Tweets Scraped: {total_tweets_all_iterations}\n
549 ")
550 #—— Sentiment analysis —————
551
552 flair_sentiment = flair.models.TextClassifier.load('en-sentiment') #flair
553 analyzer = SentimentIntensityAnalyzer() #vader
554 classifier = pipeline("sentiment-analysis",
555                        model = "finiteautomata/bertweet-base-sentiment-analysis")
556 afinn = Afinn()
557
558 def senti_score(n):
559     s = flair.data.Sentence(n)
560     flair_sentiment.predict(s)
561     total_sentiment = s.labels[0]
562     assert total_sentiment.value in ['POSITIVE', 'NEGATIVE']
563     sign = 1 if total_sentiment.value == 'POSITIVE' else -1
564     score = total_sentiment.score
565     return total_sentiment.value[0:3], round(sign * score, 3), sign
566
567 def Blob(comment):
568     try:
569         blob_Lt = TextBlob(str(comment))
570         #comment_language = detect(str(comment))
571         blob=blob_Lt.translate(from_lang='lt', to='en')
572     except:
573         blob = "Testing, delete"
574     return blob
575

```

```

576 def Textblob_fun(blob):
577     try:
578         if blob.sentiment.subjectivity > 0.3:
579             if blob.sentiment.polarity > 0:
580                 blob_values=1
581             else:
582                 blob_values=-1
583         else:
584             blob_values=0
585     except:
586         blob_values=0
587     return(blob_values)
588
589 def Vader_fun(blob):
590     try:
591         vs = analyzer.polarity_scores(blob)
592
593         if vs['compound'] >= 0.05:
594             vs_value = 1
595
596         elif vs['compound'] < -0.05:
597             vs_value = -1
598         else:
599             vs_value = 0
600     except:
601         vs_value = 0
602
603     return(vs_value)
604
605 def Transformers_fun(blob):
606     try:
607         transformer_output = classifier(str(blob[0:200]))
608         tr_result = transformer_output[0].get('label')
609         if tr_result == 'POS':
610             tr_value = 1
611         elif tr_result == 'NEG':
612             tr_value = -1
613         else:
614             tr_value = 0
615     except:
616         tr_value = 0
617
618     return(tr_value)
619
620 def Flair_fun(blob):
621     try:
622         total_sentiment, total_score, value = senti_score(str(blob))
623     except:
624         value = 0

```

```

625
626     return(value)
627
628 def Afinn_fun(blob):
629     try:
630         afinn_score = afinn.score(str(blob))
631         if afinn_score > 0:
632             afinn_value = 1
633         elif afinn_score < 0:
634             afinn_value = -1
635         else:
636             afinn_value = 0
637     except:
638         afinn_value = 0
639
640     return(afinn_value)
641
642 def sentiments(df_posts, main_folder_dr, excel_name, date):
643
644     #tracking time and fails
645     failed = 0
646     fail = 0
647     total = 0
648     start_position = 0
649     results_count = 0
650     df_with_sentiments = pd.DataFrame()
651     posts_ids = []
652     post_times = []
653     authors = []
654     posts = []
655     blobs = []
656     blob_values = [] # Textblob
657     vs_values = [] # Vader
658     tr_values = [] # Transformers
659     flair_values = [] # Flair
660     afinn_values = [] # Afinn
661
662     #Create an excel file
663     filepath = main_folder_dr+excel_name+date+'.xlsx'
664     wb = openpyxl.Workbook()
665     wb.save(filepath)
666
667     for index, row in df_posts.iterrows():
668
669         try:
670             #Post_ID
671             post_id = row['Tweet_ID']
672             posts_ids.append(post_id)
673

```

```

674     #Time
675     post_time = row[ 'Created_At' ]
676     post_times.append(str(post_time))
677
678     #Author
679     author = row[ 'Author_ID' ]
680     authors.append(str(author))
681
682     #Blob
683     comment = row[ 'Text' ]
684     posts.append(comment)
685
686     blob = Blob(comment=comment)
687     blobs.append(str(blob))
688
689     #Textblob
690     blob_value = Textblob_fun(blob)
691     blob_values.append(blob_value)
692
693     #Vader
694     vs_value = Vader_fun(blob)
695     vs_values.append(vs_value)
696
697     #Transformers
698     tr_value = Transformers_fun(blob)
699     tr_values.append(tr_value)
700
701     #Flair
702     flair_value = Flair_fun(blob)
703     flair_values.append(flair_value)
704
705     # AFINN
706     afinn_value = AFINN_fun(blob)
707     afinn_values.append(afinn_value)
708
709     #Write data to excel
710     n = 100
711     if (total%n==0 or total == len(df_posts[ 'Text' ])-failed-1) and total>0:
712         now = datetime.now().strftime("%H:%M:%S")
713         print(f'{now} {total-failed}/{total} {failed}/{n} total failed {failed}')
714
715         if total==n:
716             start_position = 0
717         else:
718             start_position = start_position + results_count
719
720         df_with_sentiments=pd.DataFrame({
721             'Tweet_ID': posts_ids,
722             'Post_It': posts,

```

```

723         'Post_en':blobs ,
724         'Created_At': post_times ,
725         'Author_ID': authors ,
726         'TextBlob values': blob_values ,
727         'Vader values': vs_values ,
728         'Transformers values': tr_values ,
729         'Flair values': flair_values ,
730         'Afinn values': afinn_values
731     })
732
733     with pd.ExcelWriter(filepath ,
734                       engine="openpyxl" ,
735                       if_sheet_exists='overlay' ,
736                       mode='a') as writer:
737         df_with_sentiments.to_excel(writer ,
738                                   startrow=start_position ,
739                                   startcol=0 ,
740                                   index=False ,
741                                   header=False)
742
743     results_count = len(blobs)
744
745     #clear variables
746     fail = 0
747     df_with_sentiments = pd.DataFrame()
748
749     posts = []
750     posts_ids = []
751     blobs = []
752     post_times = []
753     authors = []
754     blob_values = []
755     vs_values = [] # Vader
756     tr_values = [] # Transformers
757     flair_values = [] # Flair
758     afinn_values = [] # Afinn
759
760     total = total + 1
761
762     except Exception as e:
763         print(e)
764         #print(total)
765         total = total + 1
766         failed = failed + 1
767         fail = fail+1
768
769 #—— Text Preprocessing —————
770 #Remove twitter handlers
771 df_md.Text = df_md.Text.apply(lambda x:re.sub('@[\^s]+', '', str(x)))
772 #remove hashtags

```

```

772 df_md.Text = df_md.Text.apply(lambda x:re.sub(r'\B#\S+', '',x))
773 # Remove URLs
774 df_md.Text = df_md.Text.apply(lambda x:re.sub(r"http\S+", "", x))
775 # Remove all the special characters
776 df_md.Text = df_md.Text.apply(lambda x: ' '.join(re.findall(r'\w+', x)))
777 #remove all single characters
778 df_md.Text = df_md.Text.apply(lambda x:re.sub(r'\s+[a-zA-Z]\s+', '', x))
779 # Remove the '&' character
780 df_md.Text = df_md.Text.apply(lambda x: re.sub(r'&', '', x))
781 # Substituting multiple spaces with single space
782 df_md.Text = df_md.Text.apply(lambda x:re.sub(r'\s+', ' ', x, flags=re.I))
783 # Replace empty strings and spaces with NaN
784 df_md['Text'] = df_md['Text'].replace(r'^\s*$', pd.NA, regex=True)
785 # Drop rows with missing values (NaN) in the "Text" column
786 df_md.dropna(subset=['Text'], inplace=True)
787 # Reset the index if needed
788 df_md.reset_index(drop=True, inplace=True)
789
790 #—— SMI calculation —————
791
792 # convert 1 0 -1 to POS NEU NEG
793 def pos_neg_neu(df, index_name):
794     sentiments = []
795     column_name = index_name.split(' ')[0] + ' result'
796     for i in range(0, len(df)):
797         if df[index_name][i]>0:
798             sentiments.append('POS')
799         elif df[index_name][i]<0:
800             sentiments.append('NEG')
801         else:
802             sentiments.append('NEU')
803     df[column_name] = sentiments
804     return df
805
806 def senti_index(df_posts_senti, index_name='TextBlob result', lt=False):
807
808     df_senti = df_posts_senti[['Year', 'Month', index_name]]
809
810     # Count textblob_results
811     df = df_senti.groupby(['Year', 'Month', index_name]).size().reset_index(name='
Count_left')
812
813     years = df['Year'].unique()
814     months = df['Month'].unique()
815     textblob_results = df[index_name].unique()
816
817     # Create a MultiIndex DataFrame
818     result = pd.DataFrame(index=pd.MultiIndex.from_product([years, months,
textblob_results], names=['Year', 'Month', index_name]), columns=['Count_left'])

```



```

819
820 # Rename the column in the joining data
821 df_temp = df.set_index(['Year', 'Month', index_name])
822 df_temp = df_temp.rename(columns={'Count_left': 'Count'})
823
824 # Join with the original data
825 result = result.join(df_temp, how='left')
826 # Fill missing values with 0
827 result.fillna(0, inplace=True)
828
829 result = result.drop('Count_left', axis=1)
830 result = result.sort_values(by=['Year', 'Month'])
831 result.reset_index(inplace=True)
832
833 # Calculate sentiment index
834 result_neg = result[result[index_name] == 'NEG']
835 result_neg.columns = ['Year', 'Month', index_name, 'Count_neg']
836
837 result_pos = result[result[index_name] == 'POS']
838 result_pos.columns = ['Year', 'Month', index_name, 'Count_pos']
839
840 result_neu = result[result[index_name] == 'NEU']
841 result_neu.columns = ['Year', 'Month', index_name, 'Count_neu']
842
843 # Join pos with neg
844 df_pos_neg = pd.merge(result_pos, result_neg, on=['Year', 'Month'], how='inner')
845 df_pos_neg = pd.merge(df_pos_neg, result_neu, on=['Year', 'Month'], how='inner')
846
847 if lt == False:
848     column_name = index_name.split(' ')[0] + ' SI'
849 else:
850     column_name = index_name.split(' ')[0] + ' SI' + '_lt'
851 df_pos_neg[column_name] = (df_pos_neg['Count_pos'] - df_pos_neg['Count_neg']) /
852     (df_pos_neg['Count_pos'] + df_pos_neg['Count_neg'])
853 df_pos_neg["Count"] = (df_pos_neg['Count_pos'] + df_pos_neg['Count_neg'])
854 df_pos_neg["Count_all"] = (df_pos_neg['Count_pos']
855     + df_pos_neg['Count_neg']
856     + df_pos_neg['Count_neu'])
857 df_pos_neg_full = df_pos_neg
858 df_pos_neg["Year_and_Month"] = df_pos_neg['Year'].astype(str)
859     + '_'
860     + df_pos_neg['Month'].astype(str)
861
862 # Pad the 'Month' part of 'Year_and_Month'
863 # with leading zeros (if needed) to make it two digits
864 df_pos_neg['Year_and_Month'] = df_pos_neg['Year_and_Month']\
865     .str.split('_')\
866     .apply(lambda x: x[0] + '_' + x[1].zfill(2) if len(x[1]) == 1 else x[0] + '_' + x
[1])

```

```

867
868 df_pos_neg = df_pos_neg[['Year_and_Month', column_name]]
869 return pd.DataFrame(df_pos_neg), pd.DataFrame(df_pos_neg_full)
870
871 #all
872 df_TextBlob_senti, df_TextBlob_senti_full = senti_index(df_posts_senti = df_avg,
873                                                         index_name='TextBlob result')
874 df_Vader_senti, df_Vader_senti_full = senti_index(df_posts_senti = df_avg,
875                                                         index_name='Vader result')
876 df_Tr_senti, df_Transformers_senti_full = senti_index(df_posts_senti = df_avg,
877                                                         index_name='Transformers result')
878 df_Flair_senti, df_Flair_senti_full = senti_index(df_posts_senti = df_avg,
879                                                         index_name='Flair result')
880 df_Afinn_senti, df_Afinn_senti_full = senti_index(df_posts_senti = df_avg,
881                                                         index_name='Afinn result')
882
883 df_senti_list = [df_TextBlob_senti[n:m],
884                  df_Vader_senti[n:m],
885                  df_Tr_senti[n:m],
886                  df_Flair_senti[n:m],
887                  df_Afinn_senti[n:m],
888                  stat_gov_index[n:m]
889                  ]
890 df_SI = reduce(lambda left, right: pd.merge(left, right, on=['Year_and_Month'],
891                                             how='inner'), df_senti_list)
892
893
894 df_senti_list = [
895     df_TextBlob_senti_full[['Year_and_Month', 'Count']].rename(columns={"Count": "
896     Textblob_count"}),
897     df_Vader_senti_full[['Year_and_Month', 'Count']].rename(columns={"Count": "
898     Vader_count"}),
899     df_Transformers_senti_full[['Year_and_Month', 'Count']].rename(columns={"Count": "
900     Transformers_count"}),
901     df_Flair_senti_full[['Year_and_Month', 'Count']].rename(columns={"Count": "
902     Flair_count"}),
903     df_Afinn_senti_full[['Year_and_Month', 'Count', 'Count_all']].rename(columns={"
904     Count": "Afinn_count"})
905 ]
906 df_count = reduce(lambda left, right: pd.merge(left, right, on=['Year_and_Month'], how='
907     inner'), df_senti_list)
908
909 #—— SARIMAX ——
910
911 def find_stationarity_transformations(variable_series, max_diff=6):
912     """
913     Find the number of differences required to make a time series stationary.
914
915     variable_series: Pandas Series, the time series to be tested

```

```

910 max_diff: int, the maximum number of differencing allowed
911 return: A tuple containing the variable name and the number of differences required
912 """
913 for i in range(max_diff + 1):
914     # Apply differencing i times sequentially
915     differenced_series = variable_series.copy()
916     for _ in range(i):
917         differenced_series = differenced_series.diff().dropna()
918
919     # Perform the ADF test
920     if not differenced_series.empty:
921         result = adfuller(differenced_series, autolag='AIC')
922         p_value = result[1]
923         if p_value < 0.05:
924             # The series is stationary after i differencing
925             return (variable_series.name, i)
926
927 # If the series is not stationary after max_diff differencing
928 return (variable_series.name, None)
929
930 def create_differenced_dataframe(data, stationarity_transformations_df,
931 exclude_variable):
932     """
933     Create a new DataFrame with differenced variables based on the differencing order
934     specified in the stationarity_transformations_df. The column names are updated to
935     reflect the differencing order.
936
937     data: DataFrame containing the original data
938     stationarity_transformations_df: DataFrame containing variables
939         and their corresponding differencing order
940     exclude_variable: List of variables to exclude from differencing
941     return: DataFrame with differenced variables
942     """
943     data_diff = data.copy()
944
945     for col in data_diff.columns:
946         if col not in exclude_variable:
947             # Get the number of differencing required for the variable
948             diff_order = stationarity_transformations_df[
949 stationarity_transformations_df['Variable'] == col]['Differencing_Order'].iloc[0]
950
951             # Apply the differencing sequentially and update column names
952             for i in range(diff_order):
953                 data_diff[col] = data_diff[col].diff()
954             # Update the column name only for the last differencing iteration
955             if i + 1 == diff_order:
956                 new_col_name = f"{col}_diff_{diff_order}"
957                 data_diff.rename(columns={col: new_col_name}, inplace=True)

```

```

957     # Drop rows with NaN values introduced by differencing
958     data_diff = data_diff.dropna()
959
960     return data_diff
961
962 # Creating lagged variables for selected lag periods (1, 2, 3, 6, 12 months)
963 lags = [1, 2, 3, 6, 12]
964 exog_variables = data_diff_CCI.columns.drop('CCI_diff_1') # All columns except CCI
965
966 # Adding lagged variables to the dataframe
967 for lag in lags:
968     for var in exog_variables:
969         data_diff_CCI[f'{var}_lag_{lag}'] = data_diff_CCI[var].shift(lag)
970
971 # Dropping initial rows with NaN values due to lags
972 data_diff_np_lagged = data_diff_CCI.dropna()
973
974 # Calculating Pearson and Spearman correlations of lagged variables with CCI
975 pearson_corr_lagged = data_diff_np_lagged.corr(method='pearson')['CCI_diff_1']
976 spearman_corr_lagged = data_diff_np_lagged.corr(method='spearman')['CCI_diff_1']
977
978 # Filtering the results to only include the lagged variables
979 pearson_corr_lagged = pearson_corr_lagged.filter(regex='lag')
980 spearman_corr_lagged = spearman_corr_lagged.filter(regex='lag')
981
982 # CCI correlation with differentiated variables
983
984 spearman_corr_lagged_series = pd.Series(spearman_corr_lagged)
985 spearman_corr_corrected_series = pd.Series(spearman_corr_corrected)
986
987 # Extract unique variable names (excluding the '_lag_' and digits parts)
988 unique_vars = set(name.split('_lag_')[0] for name in spearman_corr_lagged.keys())
989
990 # Define the order of lags to be shown on the x-axis
991 lag_order = [0, 1, 2, 3, 6, 12]
992
993 # Find the global minimum and maximum correlation values to set a common scale
994 all_correlations = pd.concat([spearman_corr_lagged_series,
995                               spearman_corr_corrected_series])
996 vmin, vmax = all_correlations.min(), 0.35
997
998 # Define custom colormap
999 cmap = LinearSegmentedColormap.from_list("custom_cmap", ["#FFFFFF", "#78003F"])
1000
1001 # Plotting
1002 plt.figure(figsize=(10, 15))
1003 for i, var in enumerate(unique_vars):
1004     # Filter correlations for the current variable and its lags

```

```

1005     var_corr_lagged = spearman_corr_lagged_series.filter(regex=f'^{var}_lag_')
1006     var_corr_unlagged = spearman_corr_corrected_series.filter(regex=f'^{var}(?!.*_lag_)
1007     ')
1008     var_corr = pd.concat([var_corr_unlagged, var_corr_lagged])
1009
1010     # Check if the variable has any correlation data
1011     if var_corr.empty:
1012         continue
1013
1014     # Reshape for plotting
1015     var_corr_df = pd.DataFrame(var_corr, columns=[var]).T
1016
1017     # Create heatmap for the current variable
1018     ax = plt.subplot(len(unique_vars), 1, i+1)
1019     sns.heatmap(var_corr_df, annot=True, cmap=cmap, center=0, cbar=i == 0,
1020                xticklabels=lag_order, yticklabels=False,
1021                vmin=vmin,
1022                vmax=vmax,
1023                annot_kws={"size": 20})
1024     ax.set_title(f'Correlation of CCI with {var}', fontsize=17)
1025     # Increase x-axis label size
1026     ax.set_xticklabels(lag_order, fontsize=14)
1027
1028     plt.tight_layout()
1029     plt.show()
1030
1031     #Cointegration
1032     from statsmodels.tsa.vector_ar.vecm import coint_johansen
1033
1034     # Loading the differenced dataset
1035     data_diff = data_diff_CCI
1036
1037     # Selecting the stationary variables after differencing (I(1))
1038     variables_for_cointegration = ['Vader_SI_diff_1',
1039                                   'CCI_diff_1',
1040                                   'Inflation_diff_1',
1041                                   'Unemployment_without_seasonality_diff_1']
1042
1043     # Running the Johansen's cointegration test on the selected variables
1044     coint_test_result = coint_johansen(data_diff[variables_for_cointegration],
1045                                       det_order=0,
1046                                       k_ar_diff=1)
1047
1048     # Extracting the test statistics and critical values for interpretation
1049     eigenvalues = coint_test_result.eig
1050     trace_stats = coint_test_result.lr1
1051     max_eigen_stats = coint_test_result.lr2
1052     # Critical values for trace statistic
1053     critical_values_trace = coint_test_result.cvt

```

```

1053 # Critical values for max eigenvalue statistic
1054 critical_values_max_eigen = coint_test_result.cvm
1055
1056 # Prepare results for display
1057 coint_results = {
1058     "Eigenvalues": eigenvalues,
1059     "Trace Statistic": trace_stats,
1060     "Max Eigenvalue Statistic": max_eigen_stats,
1061     "Critical Values (Trace)": critical_values_trace[:, :3].tolist(),
1062     "Critical Values (Max Eigen)": critical_values_max_eigen[:, :3].tolist()
1063 }
1064
1065 # Convert the results to a DataFrame
1066 coint_results_df = pd.DataFrame(coint_results)
1067 coint_results_df
1068
1069 #STL
1070 df_senti2['Year_and_Month'] = pd.to_datetime(df_senti2['Year_and_Month'])
1071 df_senti2.set_index('Year_and_Month', inplace=True)
1072
1073 decomposition = STL(df_senti2['CCI'], period=12).fit()
1074 fig, (ax1, ax2, ax3, ax4) = plt.subplots(nrows=4, ncols=1, sharex=True, figsize=(10,8))
1075 ax1.plot(decomposition.observed, color='#E64164')
1076 ax1.set_ylabel('Observed')
1077 ax2.plot(decomposition.trend, color='#E64164')
1078 ax2.set_ylabel('Trend')
1079 ax3.plot(decomposition.seasonal, color='#E64164')
1080 ax3.set_ylabel('Seasonal')
1081 # ax4.plot(decomposition.resid)
1082 # ax4.set_ylabel('Residuals')
1083 ax4.plot(decomposition.resid.index, decomposition.resid, 'o', color='#E64164')
1084 ax4.set_ylabel('Residuals')
1085
1086 # Format the x-axis to show the dates (years)
1087 ax4.xaxis_date() # Ensure we have a date x-axis
1088 ax4.xaxis.set_major_formatter(plt.matplotlib.dates.DateFormatter("%Y")) # Show only the
    year
1089
1090 # plt.xticks(np.arange(0, 145, 12), np.arange(2018, 2024, 1))
1091 fig.autofmt_xdate()
1092 plt.tight_layout()
1093
1094 def rolling_forecast_sarimax(df, exogs, lags, p, d, q, P, D, Q, s, test_size, plot=
    False):
1095     """
1096     Perform a rolling forecast using SARIMAX
1097     and plot actual vs predicted values with dates on the x-axis.
1098
1099     df: DataFrame containing the time series

```

```

1100         and exogenous variables with 'Year_and_Month' as index.
1101     xog_variables: List of exogenous variables.
1102     lag_features: Dictionary specifying the lag for each exogenous variable.
1103     p, d, q: Non-seasonal ARIMA order parameters.
1104     P, D, Q, s: Seasonal ARIMA order parameters.
1105     test_size: Number of observations used for testing.
1106     """
1107     data = df.copy()
1108     exog_variables = exogs.copy()
1109     lag_features = lags.copy()
1110
1111     # Ensure 'Year_and_Month' is the datetime index
1112     if 'Year_and_Month' in data.columns:
1113         data['Year_and_Month'] = pd.to_datetime(data['Year_and_Month'])
1114         data.set_index('Year_and_Month', inplace=True)
1115
1116     # Create lagged features
1117     for feature, lags in lag_features.items():
1118         for lag in lags:
1119             data[f'{feature}_lag_{lag}'] = data[feature].shift(lag)
1120
1121     # Drop rows with NaN values due to lagging
1122     data = data.dropna()
1123
1124     # Update exog_variables list to include the lagged features
1125     for feature, lags in lag_features.items():
1126         exog_variables += [f'{feature}_lag_{lag}' for lag in lags]
1127
1128     # Prepare endog and exog data
1129     endog = data['CCI']
1130     exog = data[exog_variables]
1131
1132     # Split the data into training and testing sets
1133     train_endog = endog.iloc[:-test_size]
1134     train_exog = exog.iloc[:-test_size]
1135     test_endog = endog.iloc[-test_size:]
1136     test_exog = exog.iloc[-test_size:]
1137
1138     # Lists to store actual and predicted values
1139     predictions = []
1140
1141     # Rolling forecast
1142     for t in range(test_size):
1143         model = SARIMAX(train_endog,
1144                        exog=train_exog,
1145                        order=(p, d, q),
1146                        seasonal_order=(P, D, Q, s))
1147         model_fit = model.fit(dispatch=False)
1148         next_exog = test_exog.iloc[t:t+1] # Exog data for the next forecast

```

```

1149     yhat = model_fit.forecast(exog=next_exog)
1150     predictions.append(yhat.iloc[0])
1151     # Update the training set for the next iteration
1152     train_endog = train_endog.append(test_endog.iloc[t:t+1])
1153     train_exog = pd.concat([train_exog, next_exog])
1154
1155     # Calculate error metrics
1156     mse = mean_squared_error(test_endog, predictions)
1157     mae = mean_absolute_error(test_endog, predictions)
1158     rmse_val = sqrt(mse)
1159     mape_val = np.mean(np.abs((np.array(test_endog) - np.array(predictions)) / np.array
1160 (test_endog))) * 100
1161
1162     aic = model_fit.aic
1163
1164     if plot:
1165         # Print error metrics
1166         print(f"MAE: {mae}")
1167         print(f"MSE: {mse}")
1168         print(f"RMSE: {rmse_val}")
1169         print(f"MAPE: {mape_val}%")
1170
1171         # Plotting the actual values for the entire period and the forecasted values
1172         plt.figure(figsize=(14, 7))
1173         plt.plot(endog.index, endog, label='Actual', color='#414141', marker='o')
1174         forecast_index = endog.index[-test_size:]
1175         plt.plot(forecast_index,
1176                 predictions,
1177                 label='Predicted',
1178                 color='#E64164',
1179                 linestyle='—',
1180                 marker='x')
1181         plt.axvline(x=forecast_index[0],
1182                   color='#78003F',
1183                   linestyle='—',
1184                   linewidth=2,
1185                   label='Start of Test Data') # Separation line
1186         plt.title('Actual vs Predicted CCI ARMAX(1,1,1)', fontsize=32)
1187         plt.legend()
1188         plt.show()
1189
1190         # Create results DataFrame with dates
1191         results_df = pd.DataFrame({
1192             'Actual': test_endog,
1193             'Predictions': predictions
1194         }, index=test_endog.index)
1195         print(results_df)
1196
1197     return aic, mae

```



```

1197
1198 def perform_grid_search(data, exog_variables, lag_features, p_range, d, q_range,
1199                          P_range, D, Q_range, s, test_size):
1200     """
1201     Perform a grid search for SARIMAX model parameters.
1202
1203     data: DataFrame with time series data
1204     exog_variables: List of exogenous variables
1205     lag_features: Dictionary with lags for each exogenous variable
1206     p_range: Range of values for AR parameter
1207     d: Differencing parameter
1208     q_range: Range of values for MA parameter
1209     P_range: Range of values for seasonal AR parameter
1210     D: Seasonal differencing parameter
1211     Q_range: Range of values for seasonal MA parameter
1212     s: Seasonal period
1213     test_size: Size of the test dataset
1214     return: DataFrame with grid search results
1215     """
1216     results_df = pd.DataFrame(columns=['p', 'd', 'q', 'P', 'D', 'Q', 'AIC', 'MAE'])
1217
1218     for p in p_range:
1219         print('-')
1220         print(p)
1221         for q in q_range:
1222             print(q)
1223             for P in P_range:
1224                 print(P)
1225                 for Q in Q_range:
1226                     print(Q)
1227                     exog_vars = exog_variables.copy()
1228                     aic, mae = rolling_forecast_sarimax(data.copy(),
1229                                                         exog_vars,
1230                                                         lag_features,
1231                                                         p, d, q, P, D, Q, s,
1232                                                         test_size)
1233
1234                     results_df = results_df.append({
1235                         'p': p, 'd': d, 'q': q, 'P': P, 'D': D, 'Q': Q,
1236                         'AIC': aic,
1237                         'MAE': mae
1238                     }, ignore_index=True)
1239
1240     # Sort the results
1241     results_df.sort_values(by=['MAE', 'AIC'], inplace=True)
1242     return results_df
1243
1244 p_range = q_range = P_range = Q_range = range(0, 7)
1245 d = 1

```

```

1245 D = 0
1246 s = 12
1247 test_size = 7
1248 exog_variables = ['Vader_SI_diff_1',
1249                  'Inflation_diff_1',
1250                  'Unemployment_without_seasonality_diff_1']
1251 lag_features = {'Unemployment_without_seasonality_diff_1': [6],
1252                'Inflation_diff_1': [3]}
1253
1254 grid_search_results = perform_grid_search(data_diff_nCCI.copy(), exog_variables,
1255                                         lag_features, p_range, d, q_range, P_range, D, Q_range, s, test_size)
1256
1257 # Print the top 5 model combinations
1258 print(grid_search_results.head())
1259
1260 #Forecast with the best parameters
1261 exog_variables = ['Vader_SI_diff_1',
1262                  'Inflation_diff_1',
1263                  'Unemployment_without_seasonality_diff_1']
1264 lag_features = {'Unemployment_without_seasonality_diff_1': [6],
1265                'Inflation_diff_1': [3]}
1266
1267 p, d, q, P, D, Q, s = 3, 1, 5, 0, 0, 3, 12 # Example SARIMA order parameters
1268
1269 test_size = 7 # Example test size
1270 rolling_forecast_sarimax(data_diff_nCCI.copy(),
1271                          exog_variables,
1272                          lag_features,
1273                          p, d, q, P, D, Q, s,
1274                          test_size, plot = True)
1275
1276 #Decomposition
1277 def rolling_forecast_sarimax_with_coefficients(df, exogs, lags, p, d, q, P, D, Q, s,
1278                                               test_size, plot=False):
1279     """
1280     Perform a rolling forecast using SARIMAX, saving the model coefficients at each
1281     step.
1282     """
1283     data = df.copy()
1284     exog_variables = exogs.copy()
1285     lag_features = lags.copy()
1286
1287     # Ensure 'Year_and_Month' is the datetime index
1288     if 'Year_and_Month' in data.columns:
1289         data['Year_and_Month'] = pd.to_datetime(data['Year_and_Month'])
1290         data.set_index('Year_and_Month', inplace=True)
1291
1292     # Create lagged features
1293     for feature, lags in lag_features.items():

```

```

1291     for lag in lags:
1292         data[f'{feature}_lag_{lag}'] = data[feature].shift(lag)
1293
1294 # Drop rows with NaN values due to lagging
1295 data = data.dropna()
1296
1297 # Update exog_variables list to include the lagged features
1298 for feature, lags in lag_features.items():
1299     exog_variables += [f'{feature}_lag_{lag}' for lag in lags]
1300
1301 # Prepare endog and exog data
1302 endog = data['CCI']
1303 exog = data[exog_variables]
1304
1305 # Split the data into training and testing sets
1306 train_endog = endog.iloc[:-test_size]
1307 train_exog = exog.iloc[:-test_size]
1308 test_endog = endog.iloc[-test_size:]
1309 test_exog = exog.iloc[-test_size:]
1310
1311 # Lists to store actual and predicted values, and coefficients
1312 predictions = []
1313 coefficients = []
1314
1315 # Rolling forecast
1316 for t in range(test_size):
1317     model = SARIMAX(train_endog,
1318                    exog=train_exog,
1319                    order=(p, d, q),
1320                    seasonal_order=(P, D, Q, s))
1321     model_fit = model.fit(dispatch=False)
1322     next_exog = test_exog.iloc[t:t+1] # Exog data for the next forecast
1323     yhat = model_fit.forecast(exog=next_exog)
1324     predictions.append(yhat.iloc[0])
1325     coefficients.append(model_fit.params)
1326     print(model_fit.summary())
1327     model_fit.plot_diagnostics(figsize=(10,8));
1328
1329 # Assuming model_fit is your fitted SARIMAX model
1330 residuals = model_fit.resid
1331
1332 plt.tight_layout()
1333 plt.show()
1334
1335 # Update the training set for the next iteration
1336 train_endog = train_endog.append(test_endog.iloc[t:t+1])
1337 train_exog = pd.concat([train_exog, next_exog])
1338
1339 # Calculate error metrics

```

```

1340 mse = mean_squared_error(test_endog, predictions)
1341 mae = mean_absolute_error(test_endog, predictions)
1342 rmse_val = sqrt(mse)
1343 mape_val = np.mean(np.abs((np.array(test_endog) - np.array(predictions))
1344                      / np.array(test_endog))) * 100
1345 aic = model_fit.aic
1346
1347 # Create DataFrame of coefficients
1348 coefficients_df = pd.DataFrame(coefficients, index=test_endog.index)
1349
1350 if plot:
1351     # Plotting actual vs predicted values
1352     plt.figure(figsize=(12, 6))
1353     plt.plot(endog.index, endog, label='Actual', color='blue', marker='o')
1354     forecast_index = endog.index[-test_size:]
1355     plt.plot(forecast_index,
1356             predictions,
1357             label='Predicted',
1358             color='red',
1359             linestyle='—',
1360             marker='x')
1361     plt.axvline(x=forecast_index[0],
1362               color='grey',
1363               linestyle='—',
1364               linewidth=2,
1365               label='Start of Test Data')
1366     plt.title('Actual vs Predicted CCI with SARIMAX')
1367     plt.legend()
1368     plt.show()
1369
1370 # Print error metrics and results DataFrame
1371 print(f"MAE: {mae}, MSE: {mse}, RMSE: {rmse_val}, MAPE: {mape_val}%")
1372 results_df = pd.DataFrame({'Actual': test_endog,
1373                           'Predictions': predictions},
1374                           index=test_endog.index)
1375 print(results_df)
1376
1377 return aic, mae, coefficients_df
1378
1379 p, d, q, P, D, Q, s = 3, 1, 5, 0, 0, 3, 12 # Example ARIMA order parameters
1380
1381 # Example exogenous variables and lags (adjust based on your data)
1382 exog_variables = ['Vader_SI_diff_1',
1383                  'Inflation_diff_1',
1384                  'Unemployment_without_seasonality_diff_1']
1385 lag_features = {'Unemployment_without_seasonality_diff_1': [6]}
1386
1387 # Call the function
1388 aic, mae, coefficients_df = rolling_forecast_sarimax_with_coefficients(

```

```

1389     df=data_diff_nCCI, # Replace with your DataFrame name
1390     exogs=exog_variables,
1391     lags=lag_features,
1392     p=p, d=d, q=q,
1393     P=P, D=D, Q=Q,
1394     s=s,
1395     test_size=test_size,
1396     plot=True # Set to True if you want to plot the results
1397 )
1398
1399 #----- VECM -----
1400
1401 def rolling_forecast_vecm(df, original_cci_data, exog_variables, max_lag, test_size,
1402     plot=False):
1403     """
1404     Perform a rolling forecast using VECM and plot actual vs predicted values.
1405
1406     df: DataFrame containing the time series and exogenous variables.
1407     original_cci_data: DataFrame containing the original (undifferenced) CCI data.
1408     exog_variables: List of exogenous variables.
1409     max_lag: Maximum lag order for the VECM model.
1410     test_size: Number of observations used for testing.
1411     plot: Boolean flag to plot the actual vs predicted values.
1412     """
1413     data = df.copy()
1414     original_cci = original_cci_data.copy()
1415
1416     # Ensure 'Year_and_Month' is the datetime index
1417     if 'Year_and_Month' in data.columns:
1418         data['Year_and_Month'] = pd.to_datetime(data['Year_and_Month'])
1419         data.set_index('Year_and_Month', inplace=True)
1420     if 'Year_and_Month' in original_cci.columns:
1421         original_cci['Year_and_Month'] = pd.to_datetime(original_cci['Year_and_Month'])
1422         original_cci.set_index('Year_and_Month', inplace=True)
1423
1424     # Preparing the VECM data with the specified exogenous variables and target
1425     # variable
1426     vecm_data = data[exog_variables + ['CCI_diff_1']]
1427
1428     # Rolling forecast
1429     predictions = []
1430     for t in range(test_size):
1431         # Training data excludes the test set
1432         train_data = vecm_data.iloc[:-test_size+t]
1433         # Fitting the VECM model
1434         vecm_model = VECM(train_data, k_ar_diff=max_lag, coint_rank=1)
1435         vecm_fit = vecm_model.fit()
1436         # Forecasting the next step
1437         forecast = vecm_fit.predict(steps=1)

```

```

1436     predictions.append(forecast[0, -1]) # Assuming CCI_diff_1 is the last column
1437
1438 # Calculate error metrics for integrated forecasts
1439 # Last actual undifferenced CCI value
1440 last_actual_undiff_cci = original_cci['CCI'].iloc[-test_size-1]
1441 integrated_forecasts = last_actual_undiff_cci + np.cumsum(predictions)
1442 mae_integrated = mean_absolute_error(original_cci['CCI'].iloc[-test_size:],
1443                                     integrated_forecasts)
1444
1445 actual = data['CCI_diff_1'].iloc[-test_size:]
1446 mse = mean_squared_error(original_cci['CCI'].iloc[-test_size:],
1447                           integrated_forecasts)
1448 mae = mean_absolute_error(original_cci['CCI'].iloc[-test_size:],
1449                           integrated_forecasts)
1448 #mae = mean_absolute_error(actual, integrated_forecasts)
1449 rmse_val = sqrt(mse)
1450
1451 if plot:
1452     # Integrating the differenced forecasts to obtain the level forecasts for CCI
1453     # Last actual undifferenced CCI value
1454     last_actual_undiff_cci = original_cci['CCI'].iloc[-test_size-1]
1455     integrated_forecasts = last_actual_undiff_cci + np.cumsum(predictions)
1456
1457     # Plotting the forecast results for the entire period
1458     plt.figure(figsize=(14, 7))
1459     plt.plot(original_cci['CCI'][-len(original_cci)+test_size:],
1460             label='Actual', color='#414141', marker='o')
1461     plt.plot(original_cci.index[-test_size:],
1462             integrated_forecasts,
1463             label='Predicted',
1464             color='#E64164',
1465             linestyle='—',
1466             marker='x')
1467     plt.axvline(x=original_cci.index[-test_size],
1468               color='#78003F',
1469               linestyle='—',
1470               linewidth=2,
1471               label='Start of Test Data')
1472     plt.title('Actual vs Forecasted CCI (VECM)', fontsize=20)
1473     plt.xlabel('Time Periods')
1474     plt.ylabel('CCI')
1475     plt.legend(loc='lower left')
1476     plt.show()
1477
1478     return mae, mse, rmse_val
1479
1480 # the first try
1481 exog_variables = ['Vader_SI_diff_1', 'Inflation_diff_1',
1482                  'Unemployment_without_seasonality_diff_1']

```

```

1483 max_lag = 6
1484 test_size = 7
1485 mae, mse, rmse = rolling_forecast_vecm(data_diff, original_cci_data,
1486                                       exog_variables, max_lag, test_size, plot=True)
1487 print(f"MAE: {mae}, MSE: {mse}, RMSE: {rmse}") #mae 0.14998
1488
1489 # parameters combination search
1490 def test_exog_lag_combinations(df, original_cci_data, exog_vars, max_lags, test_size):
1491     results = []
1492     # Generate all non-empty combinations of exogenous variables
1493     for i in range(1, len(exog_vars) + 1):
1494         for combo in combinations(exog_vars, i):
1495             # Test each combination with all specified lag lengths
1496             for lag in range(1, max_lags + 1):
1497                 # Convert the combo to a list
1498                 # so it can be passed to the forecasting function
1499                 exog_combo = list(combo)
1500                 mae, mse, rmse = rolling_forecast_vecm(
1501                     df, original_cci_data, exog_combo, lag, test_size, plot=False
1502                 )
1503                 # Store the results
1504                 results.append({
1505                     'combo': ''.join([var[0] for var in exog_combo]),
1506                     'mae': mae,
1507                     'max_lags': lag
1508                 })
1509
1510                 print(results[-1])
1511
1512     # Convert the results to a DataFrame and sort by MAE
1513     results_df = pd.DataFrame(results)
1514     results_df_sorted = results_df.sort_values(by='mae').reset_index(drop=True)
1515     return results_df_sorted
1516
1517 # Define the exogenous variables and parameters
1518 exog_variables = ['Vader_SI_diff_1',
1519                  'Inflation_diff_1',
1520                  'Unemployment_without_seasonality_diff_1']
1521 max_lags = 12
1522
1523 # Call the function with the datasets and parameters
1524 sorted_results = test_exog_lag_combinations(
1525     df=data_diff,
1526     original_cci_data=original_cci_data,
1527     exog_vars=exog_variables,
1528     max_lags=max_lags,
1529     test_size=test_size
1530 )
1531

```

```

1532 # the second try
1533 exog_variables = ['Unemployment_without_seasonality_diff_1']
1534 max_lag = 4
1535 mae, mse, rmse = rolling_forecast_vecm(data_diff,
1536                                       original_cci_data,
1537                                       exog_variables,
1538                                       max_lag,
1539                                       test_size,
1540                                       plot=True)
1541 print(f"MAE: {mae}, MSE: {mse}, RMSE: {rmse}") #MAE: 0.01596954
1542
1543 #----- Random Forest -----
1544
1545 def plot_predictions(df, y_pred, target_column,
1546                    test_size, prediction_label='Predicted CCI'):
1547     """
1548     Plots the actual vs predicted values. Assumes the DataFrame index is a
1549     DateTimeIndex.
1550     Parameters:
1551     df (DataFrame): The original DataFrame containing the target column.
1552     y_pred (array): The predicted values.
1553     target_column (str): The name of the target column in df.
1554     prediction_label (str): Label for the predicted values in the plot.
1555     """
1556     df2 = df.copy()
1557
1558     # Ensure 'Year_and_Month' is in a DateTime format and set as index
1559     if 'Year_and_Month' in df2.columns:
1560         df2['Year_and_Month'] = pd.to_datetime(df2['Year_and_Month'])
1561         df2.set_index('Year_and_Month', inplace=True)
1562
1563     # Determine the start date of the test data within the df
1564     actual_start_of_test = df2.index[-test_size]
1565
1566     # Assuming the last entries in df correspond to the length of y_pred
1567     dates_for_pred = df2.index[-len(y_pred):] # This should align with the test set
1568
1569     plt.figure(figsize=(14, 7))
1570     plt.plot(df2.index, df2[target_column],
1571             label='Actual ' + target_column,
1572             color='#414141', marker='o')
1573     plt.plot(dates_for_pred, y_pred,
1574             label=prediction_label,
1575             color='#E64164', linestyle='—', marker='x')
1576     plt.title('Actual vs Predicted CCI Random Forest', fontsize=16)
1577     plt.axvline(x=actual_start_of_test, color='#78003F',
1578                linestyle='—', linewidth=2, label='Start of Test Data')
1579     plt.xlabel('Date')
1580     plt.ylabel(target_column)

```



```

1580 plt.legend(loc='lower left')
1581 plt.show()
1582
1583 def plot_feature_importances(model, feature_names, importance_threshold=0.01):
1584     """
1585     Plots the feature importances of a given model,
1586     filtering out features below a certain threshold.
1587     Parameters:
1588     model: The trained model.
1589     feature_names (list): List of feature names.
1590     importance_threshold (float):
1591     The threshold for feature importance (default is 0.01, or 1%).
1592     """
1593     importances = model.feature_importances_
1594     # Filter out features with importance below the threshold
1595     filtered_indices = [i for i, imp in enumerate(importance)
1596                       if imp > importance_threshold]
1597     filtered_importances = [importances[i] for i in filtered_indices]
1598     filtered_feature_names = [feature_names[i] for i in filtered_indices]
1599     indices = np.argsort(filtered_importances)[::-1]
1600     plt.figure(figsize=(10, 4))
1601     plt.title("Feature Importances")
1602     plt.bar(range(len(filtered_feature_names)),
1603            np.array(filtered_importances)[indices],
1604            color="r", align="center")
1605     plt.xticks(range(len(filtered_feature_names)),
1606               [filtered_feature_names[i] for i in indices], rotation=90)
1607     plt.xlim([-1, len(filtered_feature_names)])
1608     plt.show()
1609
1610 def summarize_model_errors(model_errors, decimals=5):
1611     """
1612     Summarizes the errors (MAE, MSE, RMSE, MAPE)
1613     of a model into a DataFrame, with rounded values.
1614     Parameters:
1615     model_errors (dict): A dictionary containing
1616     MAE, MSE, RMSE, and MAPE for the model.
1617     decimals (int): Number of decimal places to round to.
1618     Returns:
1619     pd.DataFrame: A DataFrame summarizing the errors,
1620     with each error as a row.
1621     """
1622     error_df = pd.DataFrame(model_errors.items(),
1623                            columns=['Error Metric', 'Value'])
1624     error_df['Value'] = error_df['Value'].round(decimals)
1625     return error_df
1626
1627 def summarize_feature_importances(model, feature_names, importance_threshold=0.01):
1628     """

```

```

1629     Creates a DataFrame of the feature importances in a model,
1630     filtering out those below a certain threshold.
1631     Parameters:
1632     model: The trained model.
1633     feature_names (list): List of feature names.
1634     importance_threshold (float):
1635     The threshold for feature importance (default is 0.01, or 1%).
1636     Returns:
1637     pd.DataFrame: A DataFrame listing significant feature importances.
1638     """
1639     importances = model.feature_importances_
1640     features_df = pd.DataFrame({'Feature': feature_names, 'Importance': importances})
1641     features_df = features_df[features_df['Importance'] > importance_threshold]
1642     features_df.sort_values(by='Importance', ascending=False, inplace=True)
1643     features_df['Rank'] = range(1, len(features_df) + 1)
1644     features_df.set_index('Rank', inplace=True)
1645     features_df['Feature'] = features_df['Feature']\
1646     + ' (' + round(features_df['Importance'], 4).astype(str) + ')'
1647     return features_df[['Feature']]
1648
1649 def display_formatted_tables(error_summary, feature_importance_summary):
1650     """
1651     Displays the error summary and feature importance tables in a formatted manner.
1652     Parameters:
1653     error_summary (pd.DataFrame): DataFrame containing error metrics.
1654     feature_importance_summary (pd.DataFrame): DataFrame containing feature importances
1655     .
1656     """
1657     # Set display options for better console output
1658     pd.set_option('display.max_columns', None)
1659     pd.set_option('display.expand_frame_repr', False)
1660     pd.set_option('display.max_colwidth', None)
1661     pd.set_option('display.precision', 4)
1662
1663     print("\nFeature Importance Summary:")
1664     print(feature_importance_summary)
1665
1666     print("\nModel Error Summary:")
1667     print(error_summary)
1668
1669     # Reset display options to default
1670     pd.reset_option('display.max_columns')
1671     pd.reset_option('display.expand_frame_repr')
1672     pd.reset_option('display.max_colwidth')
1673     pd.reset_option('display.precision')
1674
1675 def display_formatted_feature_importances(feature_importances_summaries):
1676     """
1677     Displays feature importance tables

```

```

1677     for each step in a formatted manner.
1678     Parameters:
1679     feature_importances_summaries (list):
1680     List of DataFrames, each containing feature importances for a step.
1681     """
1682     # Set display options for better console output
1683     pd.set_option('display.max_columns', None)
1684     pd.set_option('display.expand_frame_repr', False)
1685     pd.set_option('display.max_colwidth', None)
1686     pd.set_option('display.precision', 4)
1687
1688     for i, summary in enumerate(feature_importances_summaries, 1):
1689         print(f"\nFeature Importance Summary at Step {i}:")
1690         print(summary)
1691
1692     # Reset display options to default
1693     pd.reset_option('display.max_columns')
1694     pd.reset_option('display.expand_frame_repr')
1695     pd.reset_option('display.max_colwidth')
1696     pd.reset_option('display.precision')
1697
1698 def collect_feature_importances(feature_importances_summaries, feature_names):
1699     """
1700     Collects feature importances from summaries
1701     and stores them in a format suitable for creating a DataFrame.
1702
1703     Parameters:
1704     feature_importances_summaries (list):
1705     List of DataFrames, each containing feature importances for a step.
1706     feature_names (list): List of feature names.
1707
1708     Returns:
1709     dict: A dictionary containing feature importances for each step,
1710     organized by feature name.
1711     """
1712     feature_importances_collected = {feature: [] for feature in feature_names}
1713
1714     for summary in feature_importances_summaries:
1715         for feature in feature_names:
1716             if feature in summary.index:
1717                 feature_importances_collected[feature].append(summary.loc[feature,
1718                                                                     'Importance'])
1719             else:
1720                 feature_importances_collected[feature].append(0)
1721
1722     return feature_importances_collected
1723
1724 def combine_feature_importances(feature_importances_collected):

```

```

1725 """
1726 Combines collected feature importances into a DataFrame.
1727
1728 Parameters:
1729 feature_importances_collected (dict): A dictionary
1730 containing feature importances for each step.
1731
1732 Returns:
1733 pd.DataFrame: A DataFrame combining feature importances across steps.
1734 """
1735 return pd.DataFrame(feature_importances_collected)
1736
1737 def tune_hyperparameters(X_train, y_train):
1738     # Define the parameter grid
1739     param_grid = {
1740         'n_estimators': [100, 200, 300],
1741         'max_depth': [10, 20, 30, None],
1742         'min_samples_split': [2, 5],
1743         'min_samples_leaf': [1, 2],
1744         'max_features': ['auto', 'sqrt']
1745     }
1746
1747     # Initialize the Random Forest model
1748     rf = RandomForestRegressor(random_state=42)
1749
1750     # Perform Randomized Search
1751     rf_random = RandomizedSearchCV(estimator=rf,
1752                                   param_distributions=param_grid,
1753                                   n_iter=10, cv=3,
1754                                   random_state=42, n_jobs=-1)
1755     rf_random.fit(X_train, y_train)
1756
1757     # Return the best parameters or the best estimator
1758     return rf_random.best_params_
1759
1760 def rolling_forecast_with_lags(df, lag_features,
1761                               n_of_lags, rolling_window,
1762                               test_size,
1763                               show_plots=True, show_tables=True):
1764     """
1765     Perform a rolling forecast with lagged features
1766     and rolling averages using Random Forest.
1767
1768     Parameters:
1769     df (pd.DataFrame): The dataset to use.
1770     lag_features (list): List of column names
1771     to create lagged and rolling average features.
1772     n_of_lags (int): Number of lag periods.
1773     rolling_window (int): Window size for rolling averages.

```

```

1774     train_ratio (float): Ratio of data to be used for training.
1775
1776     Returns:
1777     dict: A dictionary containing predictions, actuals, MAE, and RMSE.
1778     """
1779     data = df[lag_features].copy()
1780
1781     # Convert 'Year_and_Month' to datetime and set it as the index (if applicable)
1782     if 'Year_and_Month' in data.columns:
1783         data['Year_and_Month'] = pd.to_datetime(data['Year_and_Month'])
1784         data.set_index('Year_and_Month', inplace=True)
1785
1786     # Create lagged and rolling average features
1787     for feature in lag_features:
1788         for lag in range(1, n_of_lags + 1):
1789             data[f'{feature}_lag_{lag}'] = data[feature].shift(lag)
1790             data[f'{feature}_rolling_avg_{rolling_window}'] \
1791                 = data[feature].rolling(window=rolling_window).mean().shift(1)
1792
1793     # Remove NaNs created by lagging and rolling
1794     data = data.dropna()
1795     split_index = len(data) - test_size
1796     train, test = data.iloc[:split_index], data.iloc[split_index:]
1797
1798     # Get the best hyperparameters
1799     best_params = tune_hyperparameters(train.drop('CCI', axis=1), train['CCI'])
1800     print(best_params)
1801
1802     # Initialize Random Forest Regressor with best hyperparameters
1803     rf_model = RandomForestRegressor(**best_params, random_state=42)
1804
1805     # Train the model
1806     X_train = train.drop('CCI', axis=1)
1807     y_train = train['CCI']
1808     rf_model.fit(X_train, y_train)
1809
1810     # Rolling forecast
1811     predictions = []
1812     feature_importances_summaries = []
1813
1814     for i in range(len(test)):
1815         # Calculate and print rolling averages
1816         #for each variable before the prediction
1817         if show_tables:
1818             for feature in lag_features:
1819                 rolling_avg = train[feature]\
1820                     .rolling(window=rolling_window).mean().iloc[-1]
1821
1822         X_test = test.iloc[i:i+1].drop('CCI', axis=1)

```

```

1823     prediction = rf_model.predict(X_test)[0]
1824     predictions.append(prediction)
1825
1826     # Add the actual value to the training set
1827     # and re-train the model
1828     new_row = test.iloc[i]
1829     train = train.append(new_row)
1830     X_train = train.drop('CCI', axis=1)
1831     y_train = train['CCI']
1832     rf_model.fit(X_train, y_train)
1833
1834     # Collect feature importances
1835     feature_importance_summary = summarize_feature_importances(rf_model,
1836                                                                X_train.columns)
1837     feature_importances_summaries.append(feature_importance_summary)
1838
1839     # Calculate MAE, MSE, RMSE
1840     mae = mean_absolute_error(test['CCI'], predictions)
1841     mse = mean_squared_error(test['CCI'], predictions)
1842     rmse = np.sqrt(mse)
1843
1844     # Calculate errors
1845     model_errors = {
1846         'MAE': mae,
1847         'MSE': mse,
1848         'RMSE': rmse
1849     }
1850
1851     if show_plots:
1852         # Diagnostics printout
1853         print("Diagnostics:")
1854         print("Length of predictions array:", len(predictions))
1855         print("Length of test set:", len(test))
1856         print("Start date for predictions:", test.index[0])
1857         print("Start date for test set:", df.index[split_index])
1858         print("First few predicted values:", predictions[:3])
1859         print("First few actual test values:", test['CCI'].head(3).tolist())
1860
1861         # Plot predictions
1862         plot_predictions(df=df,
1863                        y_pred=predictions,
1864                        target_column='CCI',
1865                        test_size=test_size,
1866                        prediction_label='Predicted CCI')
1867
1868     # Create and display error summary
1869     error_summary = summarize_model_errors(model_errors)
1870     print("\nModel Error Summary:")
1871     print(error_summary)

```

```

1872
1873     if show_tables:
1874
1875         # Display feature importance summaries for each step
1876         display_formatted_feature_importances(feature_importances_summaries)
1877
1878     results_df = pd.DataFrame()
1879     results_df['Actual'] = df['CCI'][len(predictions):]
1880     results_df['Predictions'] = predictions
1881
1882     print(f'\n{results_df}')
1883     return {
1884         'predictions': predictions,
1885         #'actuals': actuals,
1886         'mae': mae,
1887         'mse': mse,
1888         'rmse': rmse,
1889     }
1890
1891 def rolling_forecast_with_custom_lags2(data, custom_lags, rolling_variables, test_size,
1892     show_plots=True, show_tables=True):
1893     """
1894     Perform a rolling forecast with custom lagged features
1895     and rolling average features using Random Forest.
1896     Returns:
1897     dict: A dictionary containing predictions, actuals, MAE, MSE, and RMSE.
1898     """
1899     df = data.copy()
1900
1901     # Ensure 'Year_and_Month' is in datetime format and set as index
1902     if 'Year_and_Month' in df.columns:
1903         df['Year_and_Month'] = pd.to_datetime(df['Year_and_Month'])
1904         df.set_index('Year_and_Month', inplace=True)
1905
1906     # Create custom lagged features
1907     for feature, lags in custom_lags.items():
1908         for lag in lags:
1909             if lag == 0:
1910                 # Include the current period's value if lag is 0
1911                 df[f'{feature}_lag_{lag}'] = df[feature]
1912             else:
1913                 df[f'{feature}_lag_{lag}'] = df[feature].shift(lag)
1914
1915     # Create rolling average features
1916     for feature, window in rolling_variables.items():
1917         df[f'{feature}_rolling_avg_{window}'] \
1918             = df[feature].rolling(window=window).mean().shift(1)
1919

```

```

1920 # Combine selected features
1921 lag_features = [f"{feature}_lag_{lag}" for feature, lags in custom_lags.items() for
    lag in lags]
1922 rolling_features = [f"{feature}_rolling_avg_{window}" for feature, window in
    rolling_variables.items()]
1923 selected_features = lag_features + rolling_features
1924
1925 # Keep only the selected features and the target variable 'CCI'
1926 data = df[['CCI'] + selected_features].dropna()
1927
1928 # Split dataset into training and testing sets
1929 split_index = len(data) - test_size
1930 train, test = data.iloc[:split_index], data.iloc[split_index:]
1931
1932 # Prepare training data
1933 X_train = train.drop('CCI', axis=1)
1934 y_train = train['CCI']
1935
1936 # Get the best hyperparameters
1937 best_params = tune_hyperparameters(X_train, y_train)
1938
1939 # Initialize and fit the Random Forest Regressor
1940 rf_model = RandomForestRegressor(**best_params, random_state=42)
1941 rf_model.fit(X_train, y_train)
1942
1943 # Rolling forecast
1944 predictions = []
1945 feature_importances_summaries = []
1946 for i in range(len(test)):
1947     X_test = test.iloc[i:i+1].drop('CCI', axis=1)
1948     prediction = rf_model.predict(X_test)[0]
1949     predictions.append(prediction)
1950
1951     # Add actual value to the training set and re-train
1952     new_row = test.iloc[i]
1953     train = train.append(new_row)
1954     X_train = train.drop('CCI', axis=1)
1955     y_train = train['CCI']
1956     rf_model.fit(X_train, y_train)
1957
1958     # Collect feature importance
1959     feature_importance_summary = summarize_feature_importances(rf_model, X_train.
    columns)
1960     feature_importances_summaries.append(feature_importance_summary)
1961
1962 # Calculate errors
1963 model_errors = {
1964     'MAE': mean_absolute_error(test['CCI'], predictions),
1965     'MSE': mean_squared_error(test['CCI'], predictions),

```



```

1966         'RMSE': np.sqrt(mean_squared_error(test['CCI'], predictions))
1967     }
1968
1969     # Plot predictions and display tables if requested
1970     if show_plots:
1971         plot_predictions(df, predictions, 'CCI', test_size, 'Predicted CCI')
1972         error_summary = summarize_model_errors(model_errors)
1973         print("\nModel Error Summary:")
1974         print(error_summary)
1975     if show_tables:
1976         display_formatted_feature_importances(feature_importances_summaries)
1977
1978     results_df = pd.DataFrame()
1979     results_df['Actual'] = df['CCI'][len(predictions):]
1980     results_df['Predictions'] = predictions
1981
1982     return {
1983         'predictions': predictions,
1984         'model_errors': model_errors,
1985         'results_df': results_df
1986     }
1987
1988 def create_feature_importance_subplot(data_steps):
1989     """
1990     Creates a 4x2 subplot of horizontal bar charts to display feature importance for
1991     each step.
1992     Parameters:
1993     - data_steps (dict): A dictionary with steps as keys and lists of (feature,
1994     importance) tuples as values.
1995     """
1996     # Creating the subplot
1997     fig, axes = plt.subplots(4, 2, figsize=(15, 20))
1998     axes = axes.flatten() # Flatten the axes array for easy iteration
1999
2000     # Plotting data for each step
2001     for i, (step, values) in enumerate(data_steps.items()):
2002         ax = axes[i]
2003         df = pd.DataFrame(values, columns=["Feature", "Importance"])
2004         bars = ax.barh(df["Feature"], df["Importance"], color='#78003F')
2005         ax.set_title(f'{step}')
2006         ax.set_xlabel('Importance')
2007         ax.invert_yaxis()
2008
2009     # Add text labels
2010     for index, bar in enumerate(bars):
2011         # Place the text inside the bar for the first bar, otherwise outside
2012         if index == 0: # First bar
2013             text_x_pos = bar.get_width() / 1.15 # 2 Center of the bar
2014             text_color = 'white'

```

```

2013         else:
2014             text_x_pos = bar.get_width() + 0.005 # Slightly outside the bar
2015             text_color = 'black'
2016
2017             ax.text(text_x_pos, bar.get_y() + bar.get_height() / 2,
2018                    f'{bar.get_width():.4f}', va='center', color=text_color, fontsize
2019                    =8)
2020
2021 # Hide the empty subplot if the number of steps is odd
2022 if len(data_steps) % 2 != 0:
2023     axes[-1].axis('off')
2024
2025 # Adjust layout
2026 plt.tight_layout()
2027 plt.show()
2028
2029 #call RF with all features and their lags up to 12
2030 custom_lags = {
2031     'TextBlob_SI': list(range(0, 13)),
2032     'Vader_SI': list(range(0, 13)),
2033     'Transformers_SI': list(range(0, 13)),
2034     'Flair_SI': list(range(0, 13)),
2035     'Afinn_SI': list(range(0, 13)),
2036     'CCI': list(range(1, 13)), # Assuming CCI is the target, start from lag 1
2037     'Average_wage': list(range(0, 13)),
2038     'Pension': list(range(0, 13)),
2039     'Inequality': list(range(0, 13)),
2040     'Unemployment_without_seasonality': list(range(0, 13)),
2041     'Unemployment_rate': list(range(0, 13))
2042 }
2043
2044 rolling_variables = {}
2045
2046 results = rolling_forecast_with_custom_lags2(data=df_senti2,
2047                                             custom_lags=custom_lags,
2048                                             rolling_variables=rolling_variables,
2049                                             test_size=7,
2050                                             show_plots=True,
2051                                             show_tables=True)
2052
2053 #call RF with selected features and their lags
2054 custom_lags = {
2055     'CCI': [1],
2056     'TextBlob_SI': [0, 1, 7],
2057     'Inflation': [3],
2058     'Pension': [3],
2059 }
2060

```

```

2061 rolling_variables = {}
2062
2063 results = rolling_forecast_with_custom_lags2(data=df_senti2, custom_lags=custom_lags,
2064         rolling_variables=rolling_variables, test_size=7, show_plots=True, show_tables=True)
2065 #call RF with selected features, their lags and rolling averages
2066 rolling_variables = {
2067     'Inflation': 4,
2068     'Pension': 5,
2069     'CCI': 4,
2070
2071 }
2072
2073 results = rolling_forecast_with_custom_lags2(data=df_senti2, custom_lags=custom_lags,
2074         rolling_variables=rolling_variables, test_size=7, show_plots=True, show_tables=True)
2075 #----- XGBoost -----
2076
2077 def rolling_forecast_xgboost(data, target_column, date_column, lag_features, test_size,
2078         rolling_variables, plot=True):
2079     """
2080     Performs rolling forecasting with an XGBoost model.
2081
2082     Parameters:
2083     - data: DataFrame containing the dataset.
2084     - target_column: The name of the target variable column.
2085     - date_column: The name of the date column.
2086     - lag_features: Dictionary of features with their respective lags (including 0 for
2087     current values).
2088     - test_size: Number of last values to forecast.
2089     - rolling_variables: Dictionary of features with their respective rolling window
2090     sizes.
2091
2092     Returns:
2093     - pd.DataFrame: A dataframe with actual and predicted values.
2094     """
2095     df = data.copy()
2096     df_sorted = df.sort_values(by=date_column).reset_index(drop=True)
2097
2098     # Prepare the dataset with lagged and rolling features
2099     for feature, lags in lag_features.items():
2100         for lag in lags:
2101             column_name = f'{feature}_lag_{lag}'
2102             df_sorted[column_name] = df_sorted[feature].shift(lag)
2103
2104     for feature, window in rolling_variables.items():
2105         column_name = f'{feature}_rolling_avg_{window}'
2106         df_sorted[column_name] = df_sorted[feature].rolling(window=window).mean().shift

```

(1)

```

2104
2105 # Drop rows with NaN values
2106 df_sorted.dropna(inplace=True)
2107
2108 # Define features for the model
2109 model_features = [f'{feature}_lag_{lag}' for feature in lag_features for lag in
2110 lag_features[feature]] + \
2111 [f'{feature}_rolling_avg_{rolling_variables[feature]}' for feature
2112 in rolling_variables]
2113
2114 # Splitting the data into training and testing sets
2115 train_size = len(df_sorted) - test_size
2116 results_df = pd.DataFrame(columns=[date_column, 'Actual', 'Predictions'])
2117
2118 for step in range(test_size):
2119     X_train = df_sorted.iloc[:train_size + step][model_features]
2120     y_train = df_sorted.iloc[:train_size + step][target_column]
2121
2122     # Train the XGBoost model
2123     model = xgb.XGBRegressor(objective='reg:squarederror')
2124     model.fit(X_train, y_train)
2125
2126     # Plot feature importance after each step
2127     #plot_feature_importance(model, model_features, step, target_column)
2128     if plot:
2129         # Summarize feature importance after each step
2130         summarize_feature_importance(model, model_features, step)
2131
2132     # Prepare the next test point
2133     next_test_point = df_sorted.iloc[[train_size + step]][model_features]
2134     prediction = model.predict(next_test_point)[0]
2135
2136     # Append the prediction and actual value to the results dataframe
2137     results_df = results_df.append({
2138         date_column: df_sorted.iloc[train_size + step][date_column],
2139         'Actual': df_sorted.iloc[train_size + step][target_column],
2140         'Predictions': prediction
2141     }, ignore_index=True)
2142
2143 if plot:
2144     # Calculate and print metrics
2145     mae = mean_absolute_error(results_df['Actual'], results_df['Predictions'])
2146     mse = mean_squared_error(results_df['Actual'], results_df['Predictions'])
2147     rmse = np.sqrt(mse)
2148     print(f"\nMean Absolute Error: {mae}")
2149     print(f"Mean Squared Error: {mse}")
2150     print(f"Root Mean Squared Error: {rmse}")
2151
2152 # Plotting actual vs predicted values and feature importance

```

```

2151     plot_results_and_feature_importance(df_sorted, results_df, model, date_column,
2152 target_column, train_size, model_features)
2153
2154     return results_df
2155 def plot_results_and_feature_importance(df_sorted, results_df, model, date_column,
2156 target_column, train_size, model_features):
2157     """
2158     Plots the actual vs predicted values and feature importance.
2159     """
2160     plt.figure(figsize=(12, 6))
2161     plt.plot(df_sorted[date_column], df_sorted[target_column], label='Actual', color='
2162 #414141', marker='o')
2163     plt.plot(results_df[date_column], results_df['Predictions'], label='Predicted',
2164 color='#E64164', linestyle='dashed', marker='x')
2165     plt.axvline(x=df_sorted[date_column].iloc[train_size], color='#78003F', linestyle='
2166 —', label='Start of Test Data')
2167     plt.xlabel(date_column)
2168     plt.ylabel(target_column)
2169     plt.title(f'Actual vs Predicted {target_column} - XGBoost Rolling Forecast')
2170     plt.legend()
2171     plt.show()
2172
2173     feature_importance = model.feature_importances_
2174     feature_importance_df = pd.DataFrame({'Feature': model_features, 'Importance':
2175 feature_importance})
2176     feature_importance_df = feature_importance_df[feature_importance_df['Importance'] >
2177 0.01]
2178     feature_importance_df.sort_values(by='Importance', ascending=False, inplace=True)
2179
2180     plt.figure(figsize=(10, 6))
2181     plt.barh(feature_importance_df['Feature'], feature_importance_df['Importance'])
2182     plt.xlabel('Importance')
2183     plt.ylabel('Feature')
2184     plt.title('Feature Importance in XGBoost Model')
2185     plt.gca().invert_yaxis()
2186     plt.show()
2187
2188 def plot_feature_importance(model, features, step, target_column):
2189     """
2190     Plots the feature importance for the given model.
2191     """
2192     feature_importance = model.feature_importances_
2193     feature_importance_df = pd.DataFrame({'Feature': features, 'Importance':
2194 feature_importance})
2195     feature_importance_df = feature_importance_df[feature_importance_df['Importance'] >
2196 0.01]
2197     feature_importance_df.sort_values(by='Importance', ascending=False, inplace=True)

```

```

2191 plt.figure(figsize=(14, 7))
2192 bars = plt.barh(feature_importance_df[ 'Feature' ], feature_importance_df[ 'Importance
2193     '], color='#78003F')
2194 # plt.barh(feature_importance_df[ 'Feature' ], feature_importance_df[ 'Importance' ])
2195 plt.xlabel( 'Importance' )
2196 plt.ylabel( 'Feature' )
2197 plt.title(f'Feature Importance in Step {step + 1} for {target_column}')
2198
2199 # Adding percentage text at the end of each bar
2200 # Adding the percentage labels at the end of each bar
2201 for bar in bars:
2202     plt.text(bar.get_width(), bar.get_y() + bar.get_height() / 2, f'{bar.get_width
2203         (:).2%}',
2204             va='center', ha='left')
2205 plt.gca().invert_yaxis()
2206 plt.show()
2207
2208 def summarize_feature_importance(model, model_features, step):
2209     """
2210     Prints a summary of the feature importance for the given model.
2211     """
2212     feature_importance = model.feature_importances_
2213     feature_importance_df = pd.DataFrame({'Feature': model_features, 'Importance':
2214         feature_importance})
2215     feature_importance_df = feature_importance_df[feature_importance_df[ 'Importance' ] >
2216         0.01]
2217     feature_importance_df.sort_values(by='Importance', ascending=False, inplace=True)
2218     feature_importance_df[ 'Rank' ] = range(1, len(feature_importance_df) + 1)
2219     feature_importance_df[ 'Importance' ] = feature_importance_df[ 'Importance' ].round(4)
2220     feature_importance_df.set_index( 'Rank', inplace=True)
2221
2222     print(f"\nFeature Importance Summary at Step {step + 1}:")
2223     print(feature_importance_df)
2224
2225 def create_feature_importance_subplot(data_steps):
2226     """
2227     Creates a 4x2 subplot of horizontal bar charts to display feature importance for
2228     each step.
2229
2230     Parameters:
2231     - data_steps (dict): A dictionary with steps as keys and lists of (feature,
2232         importance) tuples as values.
2233     """
2234     # Creating the subplot
2235     fig, axes = plt.subplots(4, 2, figsize=(15, 20))
2236     axes = axes.flatten() # Flatten the axes array for easy iteration
2237
2238     # Plotting data for each step

```

```

2234 for i, (step, values) in enumerate(data_steps.items()):
2235     ax = axes[i]
2236     df = pd.DataFrame(values, columns=["Feature", "Importance"])
2237     bars = ax.barh(df["Feature"], df["Importance"], color='#78003F')
2238     ax.set_title(f'{step}')
2239     ax.set_xlabel('Importance')
2240     ax.invert_yaxis()
2241
2242     # Add text labels
2243     for index, bar in enumerate(bars):
2244         # Place the text inside the bar for the first bar, otherwise outside
2245         if index == 0: # First bar
2246             text_x_pos = bar.get_width() / 1.15 # 2 Center of the bar
2247             text_color = 'white'
2248         else:
2249             text_x_pos = bar.get_width() + 0.005 # Slightly outside the bar
2250             text_color = 'black'
2251
2252         ax.text(text_x_pos, bar.get_y() + bar.get_height() / 2,
2253                f'{bar.get_width():.4f}', va='center', color=text_color, fontsize
2254                =8)
2255
2256     # Hide the empty subplot if the number of steps is odd
2257     if len(data_steps) % 2 != 0:
2258         axes[-1].axis('off')
2259
2260     # Adjust layout
2261     plt.tight_layout()
2262     plt.show()
2263
2264 def find_best_rolling_windows_xgboost(data, target_column, date_column, lag_features,
2265 test_size, min_window=2, max_window=6):
2266     rolling_combinations = list(product(range(min_window, max_window + 1), repeat=len(
2267     lag_features)))
2268     results_list = []
2269
2270     for combo in rolling_combinations:
2271         rolling_variables = {feature: window for feature, window in zip(lag_features.
2272     keys(), combo)}
2273
2274         results = rolling_forecast_xgboost(
2275             data=data,
2276             target_column=target_column,
2277             date_column=date_column,
2278             lag_features=lag_features,
2279             test_size=test_size,
2280             rolling_variables=rolling_variables,
2281             plot = False
2282         )

```

```

2279
2280     mae = mean_absolute_error(results[ 'Actual' ], results[ 'Predictions' ])
2281     print(f"Combination: {combo}, MAE: {mae}")
2282     results_list.append({
2283         'Combo': combo,
2284         'MAE': mae
2285     })
2286
2287     results_df = pd.DataFrame(results_list)
2288     return results_df.sort_values(by='MAE')
2289
2290 #call XGBoost with all features and their lags up to 12
2291 lag_features = {
2292     'TextBlob_SI': list(range(0, 13)),
2293     'Vader_SI': list(range(0, 13)),
2294     'Transformers_SI': list(range(0, 13)),
2295     'Flair_SI': list(range(0, 13)),
2296     'Afinn_SI': list(range(0, 13)),
2297     'CCI': list(range(1, 13)), # Assuming CCI is the target, start from lag 1
2298     'Average_wage': list(range(0, 13)),
2299     'Pension': list(range(0, 13)),
2300     'Inequality': list(range(0, 13)),
2301     'Inflation': list(range(0, 13)),
2302     'Unemployment_without_seasonality': list(range(0, 13)),
2303     'Unemployment_rate': list(range(0, 13))
2304 }
2305 rolling_variables = {}
2306
2307 forecast_results = rolling_forecast_xgboost(
2308     data=cci_data,
2309     target_column='CCI',
2310     date_column='Year_and_Month',
2311     lag_features=lag_features,
2312     test_size=7,
2313     rolling_variables=rolling_variables
2314 )
2315
2316 #call XGBoost with selected features and their lags
2317 lag_features = {
2318     'CCI': [1],
2319     'Average_wage': [0],
2320     'Transformers_SI': [2],
2321     'Unemployment_without_seasonality':[0, 1],
2322 }
2323
2324 rolling_variables = {}
2325
2326 forecast_results = rolling_forecast_xgboost(
2327     data=cci_data,

```



```

2328     target_column='CCI' ,
2329     date_column='Year_and_Month' ,
2330     lag_features=lag_features ,
2331     test_size=7,
2332     rolling_variables=rolling_variables
2333 )
2334 #find the next rolling windows
2335 # Example usage
2336 lag_features = {
2337     'CCI': [1],
2338     'Average_wage': [0],
2339     'Transformers_SI': [2],
2340     'Unemployment_without_seasonality': [0, 1],
2341 }
2342
2343 best_rolling_windows_df = find_best_rolling_windows_xgboost(
2344     data=cci_data ,
2345     target_column='CCI' ,
2346     date_column='Year_and_Month' ,
2347     lag_features=lag_features ,
2348     test_size=7,
2349     min_window=2,
2350     max_window=12
2351 )
2352
2353 #call XGBoost with selected features , their lags and rolling averages
2354 lag_features = {
2355     'CCI': [1],
2356     'Average_wage': [0],
2357     'Transformers_SI': [2],
2358     'Unemployment_without_seasonality': [0, 1],
2359 }
2360
2361 rolling_variables = {
2362     'CCI': 6,
2363     'Average_wage': 9,
2364     'Transformers_SI': 4,
2365     'Unemployment_without_seasonality': 4,
2366 }
2367
2368 forecast_results = rolling_forecast_xgboost(
2369     data=cci_data ,
2370     target_column='CCI' ,
2371     date_column='Year_and_Month' ,
2372     lag_features=lag_features ,
2373     test_size=7,
2374     rolling_variables=rolling_variables
2375 )

```