ERASMUS UNIVERSITY ROTTERDAM
Erasmus School of Economics

Bachelor Thesis Econometrics and Operations Research
Random forest algorithms as feature selection method in macroeconomic
forecasting

Name student: van den Heuij, Pieter Thomas
Student ID number: 511648ph
date: 03-07-2022
Supervisor: Naghi, A. A
Second assessor: Bailllon, A

The views stated in this thesis are those of the author and not necessarily those of the supervisor, second assessor, Erasmus School of Economics or Erasmus University Rotterdam. We want to acknowledge all that were involved in providing guidance while carrying out this research.

**Abstract**

A frequently used method for forecasting of non-linear time series is the use of threshold autoregressive models. However, a recent study shows that machine learning could improve macroeconomic forecasting. We will be replicating elements of this study, focusing on the key features 'non-linearity' and 'regularization'. We will be using the root mean square prediction error (RMSPE) to compare Self-exciting threshold autoregressive (SETAR) models to the random forest autoregressive (RFAR) models. With this we will be answering the question "Do non-linear machine learning algorithms provide better macroeconomic forecasts than standard non-linear models?". In order to answer the question: "Does the random forest algorithm, used as a feature selection method for regularization, provide better macroeconomic forecasts?", we will compare the random forest sparse vector autoregressive (RFSVAR) model to the least absolute shrinkage and selection operator (LASSO) SVAR model, in a Data-poor as well as a Data-rich environment. We will be extending our research to multivariate time series. We utilize a dataset of macroeconomic variables concerning the U.S. economy. After concluding our research, we established that the RFAR provides significantly better forecasts for the longer horizons than the SETAR model. We also conclude that the RFSVARH+ model provides significantly better forecasts than the standard LASSO SVAR model, provided that there is no instantaneous causality.

# Contents

# 1 Introduction

Nowadays, variations of the Threshold Autoregressive (TAR) models are widely used in forecasting for non-linear time series. As Hansen (1996) explains this popularity is due to the fact that they are relatively straightforward to specify, estimate and interpret. Especially in comparison with non-linear machine learning algorithms.

Recently, the researchers from the University of Pennsylvania and the Université du Québec à Montréal have taken the discussion on the use of machine learning in macroeconomic forecasting from: 'is it useful?' to 'how can it be useful?'. They focus on four key features of machine learning algorithms, namely non-linearity, regularization, cross-validation and alternative loss function, as specified by Goulet Coulombe et al. (2022). The researchers investigate two different environments, Data-poor and Data-rich, and then apply the different features on these environments. We will be following this approach. The benefit of exploring these features, instead of finding the best model for a certain data set and certain settings, is that these features can be applied to standard models.

The paper that composes the foundation of our research is titled: "How is Machine Learning Useful for Macroeconomic Forecasting?" written by Philippe Goulet Coulombe, Maxime Leroux, Dalibor Stevanovic, and Stéphane Surprenant. They found that the feature of non-linearity is the most important for improving macroeconomic forecasting. They also conclude that the standard factor model is best for regularization, K-fold cross validation (CV) is the best CV method, and $L_2$ is the best loss function (Goulet Coulombe et al., 2022). We will be focusing on regularization as well as non-linearity, as these are the most promising for real world applications. We explore the question whether non-linear machine learning methods are indeed an improvement on standard non-linear models. A common and proven competent standard non-linear model is the Self-exciting Threshold Autoregressive (SETAR) model, as established by Feng and Liu (2003) and Boero and Lampis (2017). We will therefore use this model to research how the SETAR model compares to the random forest autoregressive (RFAR) model, as used in the paper by Goulet Coulombe et al. (2022). Hence, our first research question is stated as follows: "Do non-linear machine learning algorithms provide better macroeconomic forecasts than standard non-linear models?". The SETAR model can be useful for macroeconomic forecasting due to the fact that the U.S. economy is in a constant cycle of expansion and regression. This leads to different regimes which are accurately captured by the SETAR model, as shown by Morley (1970). The benefits of using a more standard model, over machine learning algorithms, are the possibility of easier functional utilization and lower computation times. To facilitate the comparison we start by reproducing certain results from Goulet Coulombe et al. (2022). We will then compare these results with those of the SETAR model, the root mean square prediction error (RMSPE) will be used for this comparison. Hansen (1996) solely uses the unemployment rate to compare the SETAR model with other models. Contrary to this approach, we will extend

this to five variables closely related to the U.S. business cycle. These variables match those of Goulet Coulombe et al. (2022).

We will extend the results of Goulet Coulombe et al. (2022) from the univariate time series to multivariate time series. This accounts for possible endogeneity and as Aboagye-Sarfo et al. (2015) and Zivot and Wang (2003) show, could provide improved forecasts.

For the regularization aspect of our research we will focus on the random forest sparse vector autoregressive model (RFSVAR), as proposed by Pavlyuk (2020). The promising results found by Pavlyuk (2020) are acquired in traffic forecasting. We will apply the RFSVAR in both the Data-poor as well as the Data-rich environment for macroeconomic forecasting. We denote the difference between these environments with H+ for the Data-rich environment, e.g. RFS-VARH+ is the random forest model for the Data-rich environment. This leads us to the second research question: "Does the random forest algorithm, used as a feature selection method for regularization, provide better macroeconomic forecasts?". As for the standard SVAR model we will use the Least Absolute Shrinkage and Selection Operator (LASSO) SVAR model. We will be using this model since it provides equal forecast results to RIDGE and autoregressive model with diffusion indices (ARDI) as shown by De Mol et al. (2008). In addition to this Pavlyuk (2020) suggests that the LASSO SVAR can be conveniently compared to the RFSVAR model that we will also be using.

We found that the SETAR model is not significantly outperformed by the RFAR model for the shorter horizons. For the longer horizons we do find a significant difference between the models. We also find that using multivariate times series instead of the univariate time series proposed by Goulet Coulombe et al. (2022) significantly improves the forecast accuracy. For the random forest feature selection approach we find that the RFSVARH+ model provides significantly better forecast in comparison to the standard SVARH+ and VARDI model, granted that we can ensure that there is no instantaneous causality.

Our research design contributes to answering the question presented by Goulet Coulombe et al. (2022), this being 'how can machine learning be used in macroeconomic forecasting?'. We focus our research on applying machine learning methods on the curse of dimensionality.

With our research approach we apply a regularization method that has not been performed in macroeconomic forecasting thus far. Additionally we will be extending the research done by Goulet Coulombe et al. (2022), by extending the models from univariate time series to multivariate time series.

In the remainder of the paper we will analyze the relevant literature, this is done in section 2. We continue the paper by elaborating on the data used for the replication and extension part, this will be done in section 3. Next, in section 4, we will cover the used methods and describe our approach of obtaining our results. The acquired results will be shown in section 5. Lastly, we will conclude our paper and discuss the limitations in section 6.

# 2 Literature

In this section we will explore the available literature regarding standard forecasting models and modern machine learning methods. By doing so we are able to outline current practice regarding the use of various methods. We start by elaborating on the standard forecasting models. These models have been used in the forecasting of macroeconomic variables for a longstanding period. Next, we will review literature concerning the current use of machine learning in the macro econometric field. We will also be looking into the possibilities of applying machine learning on forecasting in macro economy.

It is relevant to mention that there are different approaches of research on this matter. One common approach of research on machine learning methods is its utilization in current practice. The alternative approach is focused on the interpretation of acquired results while using machine learning methods. Our research purpose is to determine how application of machine learning can be an useful addition in the macro econometric field of forecasting. We will therefore be focusing on literature that aim attention at different ways of utilizing machine learning, instead of the interpretation of the results. The first standard forecasting model we will explore is the autoregressive (AR) model. The model originated in 1927 as found by Harrison (1999). It is a well-known model, which is oftentimes used as a benchmark model for comparison in research on macroeconomic forecasting models. Illustrated in the studies of Maehashi and Shintani (2020) and Siami-Namini et al. (2018).

The paper of Adrian et al. (2019) emphasises the significance of non-linearity in macroeconomic forecasting. While Teräsvirta (2018) establishes that there is an increase of the use of non-linear models in macro econometrics. In particular regime switching models and vector autoregressive (VAR) models. Morley (1970) depicts the relevance of allowing recurring regime switches while forecasting in macroeconomics. They demonstrate that, compared to other non-linear models, time series models that contain regime switches are the most successful form of non-linear models. The regime-switching models are mostly applied in forecasting to detect the non-linearity that arises alongside the business cycle. These models are able to identify fluctuations in the dynamics of economic activity, specifically in various stages of the business cycle. Additionally to Morley (1970) establishing that regime-switching models are the preferred model in circumstances where a business cycle is present, the existence of these business cycle is supported by Angeletos and La'O (2013). This research confirms the existence of the business cycle and therefore emphasizes the potential added benefit of using regime switching models in forecasting.

There are multiple regime switching models that can be applied for forecasting. The two most prominent regime switching models, those being Markov switching model and SETAR models, are compared in research done by Clements and Krolzig (1998). They conclude that there is no distinguishable difference between both models. We believe that the SETAR model is the most desirable due to the extensive use of SETAR models as preferred choice of regime switching model in the papers of Feng and Liu (2003) and Hansen (1999) along with

the results of Boero and Lampis (2017), which finds improved forecasts when using the SETAR model. We therefore opt to use the SETAR model as the preferred regime switching model.

A different common approach is to use the multivariate VAR models. The popularity of VAR models is explained by Pavlyuk (2020) and Zivot and Wang (2003) shows the added value of the multivariate time series model over the more standard univariate models.

Though VAR models are widely used, Pavlyuk (2020) also addresses the downside as the time dimensions increase, this downside being the curse of dimensionality as explained by Heij et al. (2004). This issue is confirmed by Davis et al. (2012). Aforementioned issues with the VAR model make for the use of sparse VAR models (SVAR) as a possible solution. By using SVAR models, we are able to bypass the time dimension issues that are present when using standard VAR models. As a result of a continuing increase of data and high time dimensions there is a corresponding rise in the need of applying SVAR models. Therefore, we will be focusing on SVAR models since the amount of data and high time dimensions are increasing and the results of SVAR models are favourable as shown by Morley (1970). We will be specifically focusing on the LASSO SVAR model, this is the multivariate versions of the model from Goulet Coulombe et al. (2022), and the random forest SVAR model from Pavlyuk (2020). The SVAR(H+) model that we research in this paper is equivalent to the $(\beta_1, \alpha = 1), K - fold$ model as proposed by Goulet Coulombe et al. (2022) in the multivariate state. We focus on the LASSO SVAR model since Pavlyuk (2020) found that this model allows for a good comparison to our RFSVAR models. De Mol et al. (2008) finds that LASSO, RIDGE, and ARDI provide similar forecast, this supports our decision to use the LASSO SVAR model, instead of replicating the ARDI model proposed by Goulet Coulombe et al. (2022). They obtain the factors in the ARDI model by principal component analysis (PCA).

We will extend on the model of Goulet Coulombe et al. (2022) by applying it in multivariate time series to account for endogenity and on the model of Pavlyuk (2020) by applying it in macroeconomic forecasting. Aboagye-Sarfo et al. (2015) and Zivot and Wang (2003) show that the multivariate systems provide for improved forecasts, therefore we will focus on the multivariate time series and use the vector autoregressive with diffusion indices (VARDI) model instead of the ARDI model as proposed by Goulet Coulombe et al. (2022). Boivin and Ng (2006) finds that increasing the available data is not beneficial and can even be harmful for estimation with factor models. The ARDI and VARDI models are examples of feature extraction methods. An alternative approach to solving the problem of high time dimensions is using machine learning methods. There is an increase in the amount of available data that needs to be processed. The importance of non-linearity in machine learning is demonstrated by Maehashi and Shintani (2020), the research shows that the random forest approach provides good results. Therefore we will select the random forest method as a feature selection method to create SVAR models. Additional evidence of the favourable performance of random forest methods is provided by

Medeiros et al. (2021). They found that the random forest method outperforms all other machine learning methods. These finding are in line with the data found in the paper of Goulet Coulombe et al. (2022). A possible explanation for this fact could be that a random forest approach is resistant to overfitting as shown by Coulombe (2020). Therefore we will replicate the RFAR model from Goulet Coulombe et al. (2022) and apply the random forest approach to enable the feature selection method for SVAR models. We will research the random forest VAR models as proposed in Pavlyuk (2020). These models have not been used in macroeconomic forecasting thus far. With this research approach we are able to improve current practice and utilisation of machine learning.

## 3 Data

In this section we will take a closer look at the data that we used for our research and show certain key features. Our research uses the Federal Reserve Economic Data-Monthly Data (FRED-MD) data set. This data set is collected by McCracken and Ng (2015). The Federal reserve bank of St. Louis chose 134 indicators that represent the U.S. macroeconomic status. They used three principles for collecting the data, these are explained by McCracken and Ng (2015). They used existing data sets, labeled these sets as vintages, and then retrieved data for the extended data sample. Thereafter they compared the newly retrieved data with the overlapping vintage data to check for any irregularities. Ultimately, they were able to compile a data set starting in 1960M01, they will be updating this set hereafter. Their goal of assembling this data set was to reduce the need for individual researchers to collect data when performing macroeconomic analysis. Furthermore, a standardized database could facilitate in replication and comparison of results. FRED-MD wanted to provide a stable data set that is readily available and updated on a regular basis so it can be used for forecasting. Comprehensive data sets, such as FRED-MD, can be used in diffusion index forecasting. It is of great value to use in research that is focused on dimension reduction for factor creating. FRED-MD had to make some adjustments to complete their data set. They check for outliers in transformed series, to ensure good factors for factor-selection procedures. McCracken and Ng (2015) defines outliers as "an observation that deviates from the sample median by more than ten interquartile ranges". They then discard the outliers they find and treat them as a missing value. Missing values are re-balanced by applying the expectation-maximization algorithm (EM-algorithm), this algorithm will be explained in section 3.1. In addition to adjusting part of the data they also provide suggestions for data transformation to obtain stationarity. In their data set there can be a column found titled 'TCODE' that contains the transformation suggestions. The Federal reserve bank of St. Louis provides the data set online, it is available to download free of charge. Coulombe downloaded the above-described historical data set from this website https://research.stlouisfed.org/econ/mccracken/sel/, as provided by McCracken and Ng (2015). They use the data set updated until

2017M12. They then apply the aforementioned transformations to achieve stationarity. Goulet Coulombe et al. (2022) mentions multiple reasons for choosing this particular data set. They explain that the FRED-MD data set has very early data available in comparison with other accessible data sets. An additional reason they provided was the fact that the components of variables are not disaggregated. Furthermore, this data set is extensively used in macroeconomic research and can therefore be frequently found in macroeconomic literature. Goulet Coulombe et al. (2022) selected five variables from the data set to carry out their research, namely 'Industrial Production' (INDPRO), 'Unemployment Rate' (UNRATE), 'Consumer Price Index' (CPI), 'difference between 10-year Treasury Constant Maturity rate and Federal Funds rate' (SPREAD), and 'housing starts' (HOUST). They decided on these specific indicators to represent the U.S. economy. They decided on using five forecasting horizons, specifically 1 month, 3 months, 9 months, 12 months, and 24 months. There are 456 evaluation periods for each horizon. When a variable exclusively contains missing data they drop that particular variable. Goulet Coulombe et al. (2022) uses an expanding window on all models while estimating. They implement an expanding window instead of a rolling window. This benefits the more flexible models as it potentially reduces their variance, as explained by Goulet Coulombe et al. (2022). We will be following this approach as it provides the ability to compare our obtained results with those of Goulet Coulombe et al. (2022).

## 3.1 Study design

The data set we used is provided by Goulet Coulombe et al. (2022) after they applied the transformations following McCracken and Ng (2015) approach. We decided on using this data set instead of the updated version of FRED-MD, with more recent data, as this provides the possibility of adequate comparison of our obtained results. It consists of 706 periods, we will be using 250 of those periods to create the models, this is equivalent to the period 1960M01 to 1979M12. The remaining 456 periods will be used as evaluation periods, covering the period 1980M1 until 2017M12. In addition to using the same data set, we also maintain other parts of the approach used by Goulet Coulombe et al. (2022) to replicate fundamental components of their research. We will be focusing on the same five variables provided in the data set, namely: INDPRO, UNRATE, CPI, SPREAD, and HOUST, as stated in Goulet Coulombe et al. (2022). We will adopt the same five forecast windows used, to reiterate 1 month, 3 months, 9 months, 12 months, and 24 months. All forecasts are made with an expanding window instead of a rolling window.

Goulet Coulombe et al. (2022) use the same approach as McCracken and Ng (2015) suggested to account for any missing values, this being the expectation-maximization algorithm (EM-algorithm). This algorithm starts by standardizing the data set and replaces all missing values with zero. By creating factors using a principal component analysis the algorithm is able to replace these zeroes with the common component of the series. The code of this algorithm was provided to us by Goulet Coulombe et al. (2022).

Meanwhile we altered part of the study design to better fit our objective. We will not be optimizing our hyperparameters, whereas Goulet Coulombe et al. (2022) intend to re-optimize their hyperparameters every two years, the main benefit of this approach is a significant reduction in computation times.

In this paper we will use the evaluation periods to create forecasts for every horizon and compare these forecasts with the true values using the RMSPE, this will be explained further in section 4.3. This comparison will be done for every model, a list of these models is provided in Table 1 with all corresponding key features. The methodology of all models is elaborated in section 4. The AR,BIC, AR,AIC, and RFAR are the models that will be replicated from Goulet Coulombe et al. (2022).

Table 1: Table featuring model characteristics

| Model | Type | Environment | Regularization | Form |
|---|---|---|---|---|
| AR,BIC | Linear | Data-poor | - | Univariate |
| AR,AIC | Linear | Data-poor | - | Univariate |
| SETAR | Non-linear | Data-poor | - | Univariate |
| RFAR | Non-linear | Data-poor | - | Univariate |
| VAR | Linear | Data-poor | - | Multivariate |
| SVAR | Linear | Data-poor | LASSO | Multivariate |
| RFSVAR | Linear | Data-poor | Random forest | Multivariate |
| VARDI | Linear | Data-rich | DI | Multivariate |
| SVARH+ | Linear | Data-rich | LASSO | Multivariate |
| RFSVARH+ | Linear | Data-rich | Random forest | Multivariate |

# 4  Methodology

In this section we will first establish which models we will replicate from the research done by Goulet Coulombe et al. (2022), thereafter we will elaborate on the models we used for the extension part of our study. We will provide the corresponding formulas and test which will be used to determine if the models are a good fit for the data. We will be using a P-value of 0.05 for all tests.

## 4.1  Replication part

We follow the approach of Goulet Coulombe et al. (2022) by defining the model for the Data-poor environment as the autoregressive direct (AR) model. We will start by elaborating on this model. Thereafter we will inspect the random forest algorithm.

We define the AR model as:

$$y_{t+h} = c + \rho(L) * y_t + \epsilon_{t+h}, t = 1, ..., T$$

with $h$ being the forecast horizon and $\rho(L)$ the lag polynomial. The order of $\rho(L)$ is the only hyperparameter that has to be optimised in this model, this

value is chosen from the subset $p_y \in \{1, 3, 6, 12\}$. We will optimize this value before creating forecasts, however we will not be optimising this value every two years contrary to Goulet Coulombe et al. (2022). This decision is made to improve the computation times drastically. To choose the optimal lag order we use the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC), this provides us with two models: AR,AIC and AR,BIC. The AIC is defined as $\log(s_p^2) + \frac{2p}{n}$ and the BIC as $\log(s_p^2) + \frac{p*\log(n)}{n}$, as defined in Heij et al. (2004). To verify if the AR models are good fits for the data we will use the Jarque-Bera (JB) test and the Ljung-Box (LB) test on the residuals to determine if these are normal distributed or serial correlated respectively. These test are further specified in Heij et al. (2004).

### 4.1.1 Random forest

Random forest was first introduced by Breiman (2001) to address the common problems found in decision trees approaches, such as overfitting and being unstable. The random forest method is robust as found by Coulombe (2020). The idea of the random forest approach is to improve the effect of averaging, e.g. lowering the variance, by reducing correlation between each tree.

We will only apply the random forest algorithm on the Data-poor environment $Z_t$, which is defined as $Z_t = [\{y_{t-j}\}_{j=0}^{p_y}]$ by Goulet Coulombe et al. (2022). We assume that all variables in $Z_t$ are covariance stationary. We will use bootstrapping to obtain different samples from $Z_t$, for each sample we then grow a random forest tree. This is done by repeatedly, for all terminal nodes, and randomly selecting m variables, picking the best split among these variables, and splitting the terminal node in to two new nodes. This process is repeated until a certain node size is reached. For further explanation of this process we refer to Hastie et al. (2009). To obtain predictions we use the average of all trees.

We use the standard value of $m = p/3$, as defined in Hastie et al. (2009), with p being the total amount of variables in the set $Z_t$. The amount of trees that we create is 500. This is sufficient to guarantee that each variable is selected at least a few times, this ensures that the prediction is stabilized as described by Goulet Coulombe et al. (2022).

We repeat the selection process of the AR model and select the number of lags from the subset $p_y \in \{1, 3, 6, 12\}$ and $p_k \in \{1, 3, 6, 12\}$ simultaneously. We determine the best lag by comparing the mean square error (MSE) of all fitted models and selecting the lowest. We then use these amount of lags to create the forecasts and compute the RMSPE. This leads to one additional model: RFAR.

## 4.2 Extension

For the extension part we will be focusing on the SETAR, VAR, SVAR, and the RFSVAR models for the Data-poor environment. For the Data-rich environment we will analyze the VARDI, SVARH+, and RFSVARH+ model.

### 4.2.1 Data-poor extension

We start by showing the equation of the SETAR model:

$$y_{t+h} = c + \rho(L)y_t * (1 - I(q_t > s)) + \kappa(H)y_t * I(q_t > s)$$

as proposed by Boero and Lampis (2017) The SETAR model allows for non-linearity due to incorporating regime switches. In this formula we have c which is a constant, and two polynomial lag orders, $\rho(L), \kappa(H)$, both lags are chosen from the set $\{1, 3, 6, 12\}$. We use the AIC to optimize these parameters. $I(q_t > s)$ indicates an indicator function this function allows for the regime switching. The threshold variable $q_t$ is a lagged endogenous variable $y_{t-d}$ in which the value of d will be automatically chosen during the estimation of the other parameters, as defined in Franses et al. (2014). This is possible due to the fact that the chosen variables closely follow the U.S. economic cycle. The threshold value s is obtained by minimizing the residual variance, defined as: $\hat{s} = arg\min\hat{\sigma}^2(s)$ with $\hat{\sigma}$ the variance of the residuals, as denoted in Franses et al. (2014). We will test the AR model against the SETAR model with 1 or 2 regimes to test for linearity by applying the F-test, for further explanation we refer to Hansen (1999). This model is estimated with the conditional least square.

The vector autoregressive (VAR) model is a model that takes into account the possible endogeneity of explanatory variables by creating a joint model of the five main variables. This a system of equations, so it is a more generalized form of the standard AR model.

We define the VAR(p) model as:

$$Y_t = \alpha + \sum_{l=1}^{p} \Phi^l Y_{t-l} + \epsilon_t, t = 1, .., T$$

as stated in Heij et al. (2004). With $Y_t$ being the vector of variables($INDPRO_t$, $UNRATE_t, CPI_t, SPREAD_t, HOUST_t$), so the amount of variables (m) is five. $\alpha$ is a vector of constants and $\Phi$ the matrix of AR coefficients of size $m*m$. The $\epsilon_t$ denotes the error term. The subscript t is the time. To determine the parameter p, we will use the AIC, defined as $AIC(p) = \log(det(\hat{\Omega}_p)) + 2\frac{pm^2}{n}$ by Heij et al. (2004). The $\hat{\Omega}_p$ is the covariance matrix of the error terms, n is the amount of observations. We assume here that all variables are stationary, this assumption holds due to the fact that McCracken and Ng (2015) provides us with stationary data. To verify if the VAR model is the right fit for the data we will check if the residual series are white noise for each variable in $Y_t$. We will do this by using the JB and Breusch-Godfred(BG) test, these test will be done in the multivariate setting. We refer to Heij et al. (2004) for an explanation of these tests. This model will be estimated with the ordinary least square (OLS) estimator. We also research the roots of the determinant of the polynomial matrix to ensure stationarity. The polynomial matrix is defined as: $\phi(Z) = I - \phi_1 z - ... - \phi_p z^p$ as defined in Heij et al. (2004).

Standard VAR models suffer from the curse of dimensionality as explained in section 2. The effect of this curse is mostly present in the Data-rich environment.

To avoid this problem we will use Sparse vector autoregressive models (SVAR), we will use the penalized least square method LASSO to reduce the amount of parameters to be estimated. This method sets certain values of the $\Phi$ to zero by introducing a penalizing function in to the objective function, we will elaborate on this topic later.

We then formulate the SVAR(p) models as:

$$Y_t = \alpha + \sum_{i=1}^{p} S^l \Phi^i Y_{t-l} + \epsilon_t$$

provided by Pavlyuk (2020). Here, $S^{(i)}$ is a binary matrix indicating the relations in the SVAR(p) model. $\alpha$ is a vector of constants and $\Phi$ a matrix of size (m x m) containing the coefficients. $Y_t$ is as specified for the VAR(p) model. $p$ is the lag order, which will be selected from the set $\{1, 3, 6, 12\}$, by computing the MSE of the CV, for all $p$, and selecting the lowest corresponding lag.

For the LASSO method the objective function with the penalty function is defined as:

$$\min_{\phi,}(\|Y_t - \mu - \sum_{i=1}^{p} S^l \Phi^i Y_{t-l}\|_F + \lambda P(\phi))$$

as specified in Pavlyuk (2020) Here we use the $L_1$ norm for the coefficient matrix and the Frobenius norm for $\|A\|_f$. We will use a CV method to determine the optimal value of $\lambda$, this is done with K-fold method, as explained by Goulet Coulombe et al. (2022). The number of folds is equal to ten, this search is done for 100 different lambda's over a reasonable grid. Again, we will perform the JB and BG test on the residuals. We will also look at the roots of the matrix $\phi(Z)$ and determine if these are outside the unit circle to ensure stationarity, as previously described for the VAR model.

For the random forest feature selection approach we will focus on the single equation strategy, as the complete model is outside the scope of this paper. Brüggemann (2004) shows that OLS for the single equation approach is efficient, when there is no instantaneous causality, it also shows certain situations were some instantaneous causality is allowed . Thus we will select the features with random forest for each equation individually, combine these into a binary matrix, to then estimate the RFSVAR model. The random forest method provides us with the increase in MSE for each feature. This is accomplished in several steps:

- Take samples of size n to create the training sets, this is done with replacement

- Use each training set to create a regression tree, where for every node in the tree the features are randomly selected.

- Calculate the importance in every tree of each feature, the increase in MSE is the measure of importance.

- Average over all training sets to obtain the feature importance.

The selection of the features is done with a sparsity of 30%, this corresponds to selecting the top 30% of features that increase the MSE the least.

For single equation the RFSVAR model is defined as:

$$y_{i,t} = \mu_i + \sum_{l=1}^{p} \sum_{j=1}^{k} s_{i,j}^l \phi_{j,t-l}^l + \epsilon_{i,t}$$

as denoted in Pavlyuk (2020).

Here, $s_{i,j}^l = 1$ when the feature is selected. $y_t$ are the individual variables from $Y_t$. The lag order, $p$, will be determined using the AIC, with $p \in \{1, 3, 6, 12\}$. k is the amount of variables in the model, for the Data-poor environment this is equal to five.

### 4.2.2 Data-rich extension

The dimension reduction techniques are most valuable in the Data-rich environment, for the VARDI model we will use the PCA component analysis to obtain the dimension reduction, for the SVARH+ model we will use the LASSO approach and for the RFSVARH+ we will use the random forest approach. For SVARH+ and RFSVARH+ we opt to only use the untransformed full data set, where as Goulet Coulombe et al. (2022) uses three different approaches. We choose this approach because the RFSVARH+ model provides interpretable results using this approach as described in Pavlyuk (2020).

The VARDI with the factors $X_t$ obtained with PCA is defined as:

$$Y_t = \alpha + \sum_{l=1}^{p} \Phi^l Y_{t-l} + \sum_{j=1}^{f} \Omega_j^* X_{t-j} \epsilon_t, t = 1, .., T$$

as provided by Bierens (2004).

In this equation $\alpha$ is a vector of constants, $\Phi^l$ the $k*k$ matrix of coefficients for lag l, and $\Omega^f$ the $k*r$ matrix of coefficients for the factors for lag j. The values of p and f denote the maximum lags of the $Y_t$ variables and the $X_t$ respectively.

In these equations, we assume that the factors, $X_t$, are exogenous, we can test this assumption with the Granger causality test as specified in Heij et al. (2004). The difference with the ARDI model from Goulet Coulombe et al. (2022) is that we will control for the endogeneity that possibly exists. We will use the JB and LB tests once more to check if the errors are white noise. Following the VAR approach we will also control the roots of the polynomial matrix. To determine the amount of factors $K$ we will compare the MSE of the models with $K \in \{3, 6, 10\}$, we follow the approach of Bai and Ng (2002). After determining the amount of factors we will determine the values of $p$ and $f$ simultaneously by computing the AIC value for all combinations of $p \in \{1, 3, 6, 12\}$ and $k \in \{1, 3, 6, 12\}$. We again choose to only optimize these parameters before forecasting and not to update them every two years. We then create the forecasts and calculate the RMSPE.

For the RFSVARH+ model, we will use the same approach as explained in 4.2.1 for the RFSVAR model. The difference between these models is that for the RFSVAR model we only use the five main variables, whereas for the RFS-VARH+ model we use all variables available from the data set, after accounting for outliers and missing values, this accounts for all possible endogeneity. We will use the single equation strategy again, thus we will select the most important features for the five main variables individually, we will use a sparsity of 30% once again. The feature weights are calculated by the random forest algorithm, as explained by Breiman (2001), this provides us with the increase in mean square error and allows us to determine the most important features.

## 4.3   Comparing the models

In the paper of Goulet Coulombe et al. (2022) the models are compared by creating forecasts for all models and calculating the root mean square prediction error (RMSPE). They create relative RMSPE with the AR,BIC model being the base model. We follow this approach in the decision to use the relative RMSPE as the evaluation criteria. First, we optimize the parameters on the data set 1960M01 until 1979M12, we then use the remaining data (1980M01-2017M12) as a pseudo out of sample data set to compute the RMSPE.

We define the RMSPE as $RMSPE = (\frac{1}{r} \sum_{h=1}^{r} (y_{n+h} - \hat{y_{n+h}})^2)^{\frac{1}{2}}$ with $r$ being the amount of observations in the hold out sample, in our case this amounts to 456 observations. $\hat{y_{n+h}}$ is the forecast for the horizon h for the particular model. $y_{n+h}$ is the true value of the variable at time $n + h$.

We compute the relative RMSPE by dividing each value using the AR,BIC as the benchmark model.

# 5   Results

This section is divided into five subsections, one for each variable that we have analyzed. We will compare our acquired results with the results of Goulet Coulombe et al. (2022). We will not be using the results of all models but focus on the models that we have replicated. This gives better comparisons between our results and those of Goulet Coulombe et al. (2022). The analysis on the differences found between these results will be done in section 6.

## 5.1   Industrial production

We begin by showing the results that we obtained in table 2, which shows the relative RMSPE for all models. The lowest values per horizon are underlined.

We clearly see that the RFSVAR model provides good forecast for the horizons $h = 1, h = 3$, and $h = 9$. The VAR model performs best for the long horizons (h=12 and h =24). When we compare the performance of the RFAR model and that of the SETAR model, we see that the SETAR model provides better forecasts for the horizons h=1 and h=3, for all other horizons the RFAR

Table 2: Relative RMSPE for the INDPRO variable

| Model | Horizon | | | | |
|---|---|---|---|---|---|
| | 1 | 3 | 9 | 12 | 24 |
| AR,BIC | 0.00656 | 0.00672 | 0.00684 | 0.00684 | 0.00684 |
| AR,AIC | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| SETAR | 1.025 | 1.030 | 1.214 | 1.356 | 2.948 |
| VAR | 0.997 | 1.010 | 1.019 | 0.993 | 0.984 |
| SVAR | 1.032 | 1.041 | 1.010 | 1.040 | 1.019 |
| RFAR | 1.043 | 1.114 | 1.187 | 1.213 | 1.175 |
| RFSVAR | 0.965 | 0.950 | 0.982 | 1.004 | 1.006 |
| VARDI | 1.535 | 1.502 | 1.348 | 1.383 | 1.391 |
| SVARH+ | 1.061 | 1.036 | 1.018 | 1.015 | 1.017 |
| RFSVARH+ | 1.193 | 1.151 | 1.110 | 1.137 | 1.113 |

model performs better. This observation could be explained by the fact that the SETAR model has a larger uncertainty regarding which regime is applicable at the longer horizons and the possibility of unit roots occurring in the regimes. When we compare the univariate models with the multivariate models we see that VAR model provides 7% and 16% lower RMSPE than the AR,BIC model for h=12 and h=24 respectively. This could be due to the fact that the VAR model captures the endogeneity that exists between the variables, this relation benefits the forecast for longer horizons. When we perform the JB and LB test for the AR,BIC and AR,AIC we see that the residuals are normally distributed though also auto correlated, therefore the VAR model provides better forecasts. We notice that the SVAR model does not provide better forecast compared to the VAR model in the Data-poor environment, this is in line with our expectations. When we compare the models in the Data-rich environment we observe a lower relative RMSPE for the SVARH+ model than the VARDI and RFS-VARH+ models for all horizons. The performance of the VARDI model can be explained by looking at JB test, which indicates that the residuals are not normally distributed due to high kurtosis. The assumption of no instantaneous causality could explain the performance of the RFSVARH+ model, due to the large amount of variables included in the Data-rich environment the change of instantaneous causality increases, this chance is lower in the Data-poor environment. Hence, we see a significantly lower RMSPE for the RFSVAR model.

Goulet Coulombe et al. (2022) found that the univariate version of the SVAR model was the best at forecasting the INDPRO variable for the horizons h=1 and h=3. For all other horizons they found that the RFAR model has the lowest RMSPE. That would indicate that the RFAR model is the best forecasting model for long horizons. However, we find that the VAR model is the preferred model at longer horizons. This shows the importance of the multivariate models in time series forecasting.

## 5.2 Unemployment rate

When we compare the results of the INDPRO variable with the UNRATE variable we find similar results, meaning that the RFSVAR model provides good forecasting results. For the UNRATE variable it is the best forecasting model for h=1, h=3, and h=12. The RFSVARH+ model outperforms the other models for h=24. Table 3 shows the RMSPE for the UNRATE variable, again with the lowest values underlined.

Table 3: Relative RMSPE for the UNRATE variable

| Model | Horizon | | | | |
|---|---|---|---|---|---|
| | 1 | 3 | 9 | 12 | 24 |
| AR,BIC | 0.163 | 0.162 | 0.171 | 0.171 | 0.171 |
| AR,AIC | 0.995 | 0.995 | 0.995 | 1.006 | 1.000 |
| SETAR | 1.032 | 1.088 | 1.210 | 1.343 | 3.177 |
| VAR | 0.988 | 1.022 | 0.999 | 0.991 | 0.988 |
| SVAR | 1.022 | 1.030 | 0.989 | 1.015 | 1.014 |
| RFAR | 1.028 | 1.036 | 1.010 | 1.018 | 1.015 |
| RFSVAR | 0.968 | 1.009 | 0.954 | 0.990 | 0.992 |
| VARDI | 1.369 | 1.456 | 1.337 | 1.290 | 1.317 |
| SVARH+ | 1.080 | 1.101 | 1.035 | 1.020 | 1.001 |
| RFSVARH+ | 1.042 | 1.083 | 1.034 | 1.029 | 0.986 |

The RFSVARH+ model performs well for the UNRATE variable, a possible explanation is that the variables that are chosen by the random forest algorithm do not have instantaneous causality. Therefore the randomisation and the Data-rich environment give significantly lower RMSPE than the VARDI and SVARH+ models. The SETAR model provides high relative RMSPE in comparison with the AR,BIC model, this can be explained by the fact that the lower regime contains roots inside the unit circle, this effects the stationarity within this regime. Similar to the INDPRO variable we notice that the VAR model significantly outperforms the AR,BIC and AR,AIC model, i.e. the VAR model captures certain endogeneity between the five main variables.

## 5.3 Consumer price index

When we evaluate the results for the variable CPI we observe different results compared to the INDPRO and UNRATE variables. We see that the SETAR model provides the lowest RMSPE for the horizons 1 and 12. Contrary to the previous variables the forecasts for longer horizons are significantly better or similar to the AR model, this can be explained by the fact that the CPI variable has a longer cycle, as described by Morley (1970). This eliminates the problem as previously described. We again see a relatively low RMSPE for the RFSVAR model and a significantly high RMSPE for the RFSVARH+ model, this is explained by the verity of the assumption of no instantaneous causality. In line with the previous variables we see that the SVARH+ model significantly

outperforms the VARDI model in the Data-rich environment. When we analyze the residuals of the VARDI model we can conclude that these are not normally distributed, as a result the estimations are not efficient.

Table 4 shows the RMSPE for the CPI variable.

Table 4: Relative RMSPE for the CPI variable

| Model | Horizon | | | | |
|---|---|---|---|---|---|
| | 1 | 3 | 9 | 12 | 24 |
| AR,BIC | 0.00260 | 0.00280 | 0.00277 | 0.00277 | 0.00276 |
| AR,AIC | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| SETAR | 0.989 | 1.018 | 0.990 | 0.994 | 1.005 |
| VAR | 1.024 | 1.039 | 0.995 | 0.999 | 1.001 |
| SVAR | 1.062 | 0.986 | 0.998 | 0.998 | 1.001 |
| RFAR | 1.218 | 1.178 | 1.099 | 1.133 | 1.080 |
| RFSVAR | 1.016 | 1.033 | 0.986 | 0.998 | 1.000 |
| VARDI | 1.226 | 1.470 | 1.541 | 1.525 | 1.491 |
| SVARH+ | 1.061 | 0.985 | 0.997 | 0.997 | 1.000 |
| RFSVARH+ | 1.722 | 1.610 | 1.640 | 1.626 | 1.652 |

Goulet Coulombe et al. (2022) finds that for this specific variable, the RFAR model works best for forecasting for $h = 9$ and $h = 12$. They find that the AR model provides the best forecast for the long horizon, $h = 24$. This is in accordance with our results. For $h = 1$ and $h = 3$ they find that the univariate version of our SVARH+ model provides the lowest RMSPE. We find similar results with the exception of horizon h=1, as a result of the reasons previously stated.

## 5.4 Housing starts

We start by showing the relative RMSPE for the variable HOUST in table 5. The lowest values are again underlined.

Here we see a highly significant result, the RFSVARH+ model provides the lowest RMSPE for the horizons $h = 9, h = 12$, and $h = 24$. We do not find instantaneous causality within the variables that are chosen with the random forest algorithm. The benefit of this comes to light for the longer horizons. For h=9 we find that the RFSVAR models provides the best forecasts. The significant lower RMSPE of the VAR model in comparison to the AR,BIC model is evidence of the fact that the five main variables are endogenous. When we examine the results of the VARDI model for the HOUST variable we find that the relative RMSPE is up to three times as high as that of the AR,BIC model. For the ARDI model we assume that the factors are exogenous, as shown by the Granger causality test this assumption does not hold and the results are therefore not efficient. For the RFSVARH+ model we find evidence of instantaneous causality.

Table 5: Relative RMSPE for the HOUST variable

| Model | Horizon | | | | |
|---|---|---|---|---|---|
| | 1 | 3 | 9 | 12 | 24 |
| AR,BIC | 0.0756 | 0.105 | 0.184 | 0.212 | 0.294 |
| AR,AIC | 0.997 | 0.990 | 1.021 | 1.028 | 1.023 |
| SETAR | 1.033 | 1.053 | 1.153 | 1.156 | 1.064 |
| VAR | 1.019 | 0.834 | 0.667 | 0.635 | 0.545 |
| SVAR | 1.128 | 1.166 | 1.061 | 1.148 | 1.095 |
| RFAR | 1.269 | 1.178 | 1.072 | 1.056 | 1.117 |
| RFSVAR | 1.001 | 0.792 | 0.659 | 0.707 | 0.619 |
| VARDI | 3.034 | 2.225 | 1.397 | 1.302 | 1.290 |
| SVARH+ | 5.074 | 3.514 | 1.980 | 1.728 | 1.278 |
| RFSVARH+ | 1.248 | 0.902 | 0.516 | 0.445 | 0.332 |

## 5.5 Spread

Lastly, we inspect the results of the spread between the 10-year maturity rate and the rate of the Federal Funds. Table 6 shows the relative RMSPE for the SPREAD variable.

Table 6: Relative RMSPE for the SPREAD variable

| Model | Horizon | | | | |
|---|---|---|---|---|---|
| | 1 | 3 | 9 | 12 | 24 |
| AR,BIC | 0.541 | 1.036 | 1.308 | 1.473 | 1.692 |
| AR,AIC | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| SETAR | 1.172 | 1.169 | 1.204 | 1.098 | 1.145 |
| VAR | 0.989 | 0.611 | 0.742 | 0.773 | 0.890 |
| SVAR | 2.196 | 1.157 | 0.815 | 0.868 | 1.004 |
| RFAR | 1.654 | 1.085 | 1.151 | 1.093 | 1.242 |
| RFSVAR | 1.043 | 0.612 | 0.640 | 0.728 | 0.829 |
| VARDI | 3.122 | 1.683 | 1.243 | 1.134 | 0.993 |
| SVARH+ | 7.585 | 3.684 | 2.380 | 1.931 | 1.298 |
| RFSVARH+ | 1.636 | 0.855 | 0.695 | 0.605 | 0.533 |

We find that the multivariate VAR model is the best forecasting model for the horizons $h = 1$ and $h = 3$, this is due to endogeneity of the variables. For the longer horizons($h = 12, h = 24$ we see that the RFSVARH+ provides the lowest RMSPE, we find no evidence of instantaneous causality which explains the performance of the RFSVARH+ model. For h=9 the RFSVAR model has the lowest RMSPE. The SVARH+ model is significantly outperformed by the AR,BIC model this can be explained by the fact that the SVARH+ model is not scale invariant. When we compare the RFAR model and the SETAR model we see no significant difference in their ability to forecast the SPREAD variable.

# 6    Conclusion

We conclude this paper by outlining our main research findings and discussing the differences between our research and that of Goulet Coulombe et al. (2022). We will summarize our main findings and answer our research questions. Subsequently we will address the limitations we encountered while conducting our research and provide further research recommendations.

The research executed by Goulet Coulombe et al. (2022) encouraged us to elaborate on the use of machine learning in macroeconomic forecasting. After they found that machine learning, especially the features non-linearity and regularization, could be crucial for improving macroeconomic forecasting we decided to extend their research by applying these approaches to a multivariate time series. We also replicated part of their study approach to be able to compare our results.

The replication part of our research consisted of creating the AR and RFAR and using the RMSPE to compare our SETAR model with these models as described by Goulet Coulombe et al. (2022). This comparison was used to define their ability of macro-economic forecasting in non-linear settings. We also compared the RFSVAR to the standard LASSO SVAR model and to the VARDI model to determine the performance of random forest algorithms in feature selection.

With the described approach we intended to answer two research questions. We started this paper by specifying our first research question: "Do non-linear machine learning algorithms provide better macroeconomic forecasts than standard non-linear models?". We decided on this subject matter as we believed that the researchers of the replicated paper did not sufficiently explore standard non-linear models. We found that the RFAR models provide an overall more accurate forecast for longer horizons, this is in line with the results of Goulet Coulombe et al. (2022). The SETAR and RFAR model do not significantly differ for shorter horizons. The answer to our research question is therefore that the RFAR algorithm does provide better macroeconomic forecasts over the standard SETAR model.

Apart from the non-linearity aspect of machine learning methods we also explored the regularization feature of machine learning. This led to our second research question, being: ""Does the random forest algorithm, used as a feature selection method for regularization, provide better macro-economic forecasts?". After concluding our research we found that RFSVARH+ outperforms both SVAR and VARDI as a feature selecting method for regularization, if we can control if there is any instantaneous causality or we are in the setting described by Brüggemann (2004).

The purpose of our research is to create a better understanding of the potential benefits that machine learning could provide in macroeconomic forecasting. Since the models which we researched have not been applied in macroeconomic forecasting thus far, our research contributes to improving the use of machine learning in current practice.

The main difference between the overall values of the RMSPE and the dif-

ferent results between our research and that of Goulet Coulombe et al. (2022) is due to two main factors. The first one is our decision to not optimise the parameter selection every two years. The other part of the difference can be explained by looking at the different transformations, we chose to use the data with the transformations of McCracken and Ng (2015) were we believe that Goulet Coulombe et al. (2022) opted to standardize and transform this data set, they did not specifically clarify this in their paper.

The limitations we encountered while performing our research mainly derived from the fact that Goulet Coulombe et al. (2022) did not adequately disclose their research approach. By omitting access to their code for the non-linear models, we experienced difficulties in reproducing their results. In addition to this we decided to deviate from the optimization method used by Goulet Coulombe et al. (2022), we did not optimize the parameters every two years since this provides lower computation times. We suggest that future research is focused on the interpretability of the machine learning methods and other multiple regime models. We suggest the Smooth Transition AR (STAR) model or the Markov-Switching (MSW) model. We also want to propose further research to determine the forecasting performance of the models in the NBER recession periods.

# References

Aboagye-Sarfo, P., Mai, Q., Sanfilippo, F. M., Preen, D. B., Stewart, L. M., and Fatovich, D. M. (2015). A comparison of multivariate and univariate time series approaches to modelling and forecasting emergency department demand in western australia. *Journal of Biomedical Informatics*, 57:62–73.

Adrian, T., Boyarchenko, N., and Giannone, D. (2019). Vulnerable Growth. *American Economic Review*, 109(4):1263–1289.

Angeletos, G.-M. and La'O, J. (2013). Sentiments. *Econometrica*, 81(2):739–779.

Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221.

Bierens, H. J. (2004). Var models with exogenous variables.

Boero, G. and Lampis, F. (2017). The forecasting performance of setar models: An empirical application. *Bulletin of Economic Research*, 69(3):216–228.

Boivin, J. and Ng, S. (2006). Are more data always better for factor analysis? *Journal of Econometrics*, 132(1):169–194. Common Features.

Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.

Brüggemann, R. (2004). *Model reduction methods for vector autoregressive processes*. Springer.

Clements, M. P. and Krolzig, H.-M. (1998). A comparison of the forecast performance of markov-switching and threshold autoregressive models of us gnp. *The Econometrics Journal*, 1(1):C47–C75.

Coulombe, P. G. (2020). To bag is to prune.

Davis, R. A., Zang, P., and Zheng, T. (2012). Sparse vector autoregressive modeling.

De Mol, C., Giannone, D., and Reichlin, L. (2008). Forecasting using a large number of predictors: Is bayesian shrinkage a valid alternative to principal components? *Journal of Econometrics*, 146(2):318–328. Honoring the research contributions of Charles R. Nelson.

Feng, H. and Liu, J. (2003). A setar model for canadian gdp: non-linearities and forecast comparisons. *Applied Economics*, 35(18):1957–1964.

Franses, P., van Dijk, D., and Opschoor, A. (2014). *Time Series Models for Business and Economic Forecasting*. Cambridge University Press.

Goulet Coulombe, P., Leroux, M., Stevanovic, D., and Surprenant, S. (2022). How is machine learning useful for macroeconomic forecasting? *Journal of Applied Econometrics*.

Hansen, B. (1999). Testing for linearity. *Journal of Economic Surveys*, 13(5):551–576.

Hansen, B. E. (1996). Estimation of TAR Models. Boston College Working Papers in Economics 325., Boston College Department of Economics.

Harrison, P. (1999). *Southern Economic Journal*, 66(1):200–202.

Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*. Springer, 2 edition.

Heij, C., de Boer, P., Franses, H., Kloek, T., and van Dijk, H. (2004). *Econometric methods with applications in business and Economics*. Oxford University Press.

Maehashi, K. and Shintani, M. (2020). Macroeconomic forecasting using factor models and machine learning: an application to japan. *Journal of the Japanese and International Economies*, 58:101104.

McCracken, M. and Ng, S. (2015). Fred-md: A monthly database for macroeconomic research.

Medeiros, M. C., Vasconcelos, G. F. R., Álvaro Veiga, and Zilberman, E. (2021). Forecasting inflation in a data-rich environment: The benefits of machine learning methods. *Journal of Business & Economic Statistics*, 39(1):98–119.

Morley, J. (1970). *Macroeconomics, Non-linear Time Series in*.

Pavlyuk, D. (2020). Random forest variable selection for sparse vector autoregressive models. In Valenzuela, O., Rojas, F., Herrera, L. J., Pomares, H., and Rojas, I., editors, *Theory and Applications of Time Series Analysis*, pages 3–17, Cham. Springer International Publishing.

Siami-Namini, S., Tavakoli, N., and Siami Namin, A. (2018). A comparison of arima and lstm in forecasting time series. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 1394–1401.

Teräsvirta, T. (2018). Nonlinear models in macroeconometrics.

Zivot, E. and Wang, J. (2003). *Vector Autoregressive Models for Multivariate Time Series*, pages 369–413. Springer New York, New York, NY.

# 7  Appendix

ReadMeFile

You first download this zipfile to your desktop. You than have to change the path in each file to download the FRED-MD_JAN-2018 Data set and set the right workdrive.

The provide zipcode file contains 12 files, 1 excel file with the data and 11 Rstudio files. The excel file, EM_sw.R file and the factor.R file are provided by Stephane Surprenant  Creation: 16/11/2017 (Goulet Coulombe et al., 2022), the EM_sw.R function creates factors using principal component analysis and replace missing data with an expectation-maximization algorithm. The factor.R file is a function which is used in the EM_sw.R file to help create the factors. Stephane Surprenant also provided use with one line of code to easily drop all the variables with only missing values prior to applying the EM algorithm.

For the results in tables 2,3,4, 5, and 6 you will have to run the following files: ARModels.R, SETAR.R, VAR.R, SVAR.R, RFAR.R, VARDI.R, RFSVAR.R, SVARH+.R, RFSVARH+.R. This has to be done for each horizon (1,3,9,12, and 24). This gives the RMSPE errors for all models and each horizon, to obtain the relative RMSPE you will have to manually calculate these. The ARModels.R file gives the RMSPE for the AR,BIC and AR,AIC model and the JB and LB test for the AR models. The SETAR.R model gives the RMSPE for the SETAR model and the linearity test. The VAR.R file provides you with the RMSPE, JB test, LB test, and the roots for the VAR model. The SVAR.R file provides you with the same informaion as the VAR.R file but than for the SVAR model. The RFAR.R file calculates the RMSPE for the RFAR model. The VARDI.R file gives the RMSPE of the VARDI model and the Granger causality test. The RFSVAR.R file gives the RMSPE for the RFSVAR model as well as the roots. The SVARH+ model calculates the RMSPE of the SVARH+ model in the Data-rich environment. The last file, RFSVARH+.R calculates the RMSPE of the RFSVARH+ model.

We emphasise the fact that all files have to be run for each horizon to obtain all RMSPE. All needed packages are provided at the top of each file.