

# Deep Vector Autoregression for Macroeconomic Data

Marc Agustí (marc.agusti@barcelonagse.eu)

Patrick Altmeyer (patrick.altmeyer@barcelonagse.eu)

Ignacio Vidal-Quadras Costa (ignacio.vidalquadrascosta@barcelonagse.eu)

junio, 2021

## Abstract

Vector autoregression models have been the traditional technique in the last decades when it comes to the forecasting of time series data. In particular, the VAR framework has demonstrated its outperformance in precision and accuracy of predicting the next lags of time series, hence becoming a useful tool for policy makers, which rely on this methodology to construct forecasts or even for the analysis of impulse response functions. However, with the recent advancement in computational power of computers, and more importantly, the development of more advanced machine learning algorithms and approaches such as deep learning, new algorithms are developed to analyze and forecast time series data. This paper aims to contribute to the time series literature by using the Long Short-Term Memory (LSTM) to see whether this new advanced methodology is superior to the usual VAR framework. By fitting each regression of the VAR with a neural network instead of with a simple OLS regression, we are able to outperform both in-sample and out-of-sample the predictions of the usual VAR for variables such as GDP, FED funds rate and inflation, ...

**Keywords**— Vector Autoregression, Deep Learning, Forecasting, Neural Networks, Macroeconomic Timeseries

# 1 Introduction

As stated by the European Central Bank, the monetary transmission mechanism is the process through which monetary policy decisions affect the economy in general and the price level in particular. However, the uncertainty of this transmission is huge, given that this transmission is characterized by long, variable and uncertain time lags. Hence, it is difficult to predict how changes in monetary policy actions affect economic real variables. Therefore, it is of foremost importance for policy-makers to come up with tools that allow them to predict these effects.

With this in mind, a lot of research on the forecasting of time series has been developed to assess the effect of current policy decisions on future economic variables. Thanks to this, over the last decades policy makers have had more information when taking decisions. This information can be of the form of point estimates and interval forecasts. To come up with these estimates, several methodologies have been largely applied in the time series forecasting literature.

Up until now, the most common methodology to come up with these estimates is the so-called Vector Autoregression (VAR). This framework, which belongs to the traditional class of econometric forecasting techniques, has proved to provide policy-makers with fairly good and consistent point and interval estimates. It has therefore been used in the monetary policy divisions of institutions such as the European Central Bank or the International Monetary Fund.

Simultaneously, with the recent advancement in computational power of computers, and more importantly, the development of more advanced machine learning algorithms and approaches such as deep learning, new algorithms are developed to analyze and forecast time series data. Whereas the good performance of techniques such as VAR is well-known, it is still uncertain whether deep learning algorithms can actually improve these forecasts.

To this end, this paper aims to provide with a new and ground-breaking methodology that combines the VAR equation-by-equation structure with the deep neural network architecture to capture the potential non-linear relationships in the data generating process. Thereupon, we construct what we called a Deep-VAR framework to generate point and interval forecasts. The ultimate objective of this empirical exercise is to provide the existing literature with a method that outperforms the most conventional methods, yet keeping interpretability in the underlying mechanism.

To the best of our knowledge, this is the first paper to fit a deep neural network for each equation of the VAR structure. Although a lot of authors have used deep learning to forecast time series, all of them fit the whole system in one big neural network. We believe that by doing it equation-by-equation not only do we maintain interpretability of the model but also we get better forecasting results.

We find that the Deep-VAR methodology outperforms the traditional VAR framework both in-sample and out-of-sample. When it comes to forecasts, we also obtain a lower RMSE.

Yet, policy makers are not just interested in the forecasting accuracy of the model but they also have an interest in the inference side. For instance, central banks need to know if rates granger cause one variable or not. Another aspect policy makers and researchers are interested in is how the variables of the system evolve through time once a shock takes place. This information is recovered using Impulse Response Functions (IRFs).

The linear additive relationship of linear models allows the model to observe which variable granger cause another and what are the IRFs of the variables in the system. Unfortunately, in the case of

nonlinear models, its nature makes it impossible to recover these insights because the outputs and the inputs do not have a linear additive structure.

Therefore, on the one hand, the nonlinear structure of NN helps us for forecasting in the case that there are nonlinear relationships present in the series, but on the other hand we lose the interpretability of the model, making it impossible to recover IRFs or to even to assess if one variable granger causes another. This problem is also known as the black box problem. This is because the data fed into the input layer passes through the succeeding layers, getting multiplied, added together and transformed in complex and different ways, until it finally arrives, radically transformed, to the output layer. Therefore, it is impossible to assess what happened with one input and how it affected the output.

The rest of the paper is structured as follows: in section 2 we present a literature review of prior research on the methodologies used to provide forecasts and on the monetary transmission mechanism. Section 3 is a detailed description of the data we use in our empirical exercise. In section 4 we present the traditional VAR methodology and we also present our Deep-VAR model following by our empirical findings in section 5. Finally, in section 6 we finish with concluding remarks.

## 2 Literature review

There is a large agreement among economists on the fact that monetary policy has a short-term influence on the economic activity. Friedman and Schwartz (2008) found that monetary policy actions are followed by movements in real output that may last for two years or more (Romer and Romer (1989); Bernanke (1990)). However, what are the forces that trigger these effects is of interest for most economists, in particular, economists aim to understand the monetary transmission mechanism. If monetary policy affects the real economy, what is the transmission mechanism by which these effects occur? This is one of the questions which is among the most important and controversial in macroeconomics.

In the aftermath of the oil price shock in the 1970's, interest was raised in understanding business cycles. To do that, most economists made use of large-scale macroeconomic models, which was criticized by Lucas Jr (1976), stating that the assumption of invariant behavioral equations was inconsistent with the dynamic maximizing behavior. Hence, New Classical economists started making use of the so-called market clearing models of economic fluctuations. With the goal of really taking into account productivity shocks, Real Business Cycle models were developed (Kydland and Prescott (1982)).

After the failure of the large-scale macroeconomic models when trying to predict business cycles, the economic profession tried to solve this using structural vector autoregression (VAR) models to analyze business cycles, which were useful to capture the impact of policy-actions. Sims and others (1986) suggested that VARs were useful to evaluate macroeconomic models. One of the advantage of VARs is that they are not a large and complicated structure, and hence are easily interpretable, that is, do not suffer from the “black box” problem.

In the last decades the use of VAR's in order to do time series forecasting has been quite extensive. Actually, a lot of different models have been proposed with the intention to model and predict time series data. When it comes to the VAR framework, the different factors in the projected VAR

models are difficult to understand, and that is why researchers rely heavily on impulse response functions (IRF) (Enders (2008)).

As for now, the models that we have seen are traditional econometric based models that are not able to capture nonlinear relationships in the data, which might be sometimes a limitation. In the case of economic time series, specifically gdp, inflation and so on, nonlinear trends are likely to appear and be present in the essence of the data generating process as shown by Brock et al. (1991), specially as a reponse to large, unexpected economic fluctuations.

In the past years, authors have therefore started using nonlinear techniques for forecasting. Machine Learning has contributed a lot to this field. The most popular machine learning techniques which do not assume a linear relationship between inputs and outputs are K-Nearest Neighbors (first introduced by Fix and Hodges (1951)), Support Vector Machines (mostly developed by Cortes and Vapnik (1995)), Random Forests (first introduced in 1995 by Ho (1995)) and Neural Networks (NN) (first proposed in 1943 by McCulloch and Pitts (1990)).

With these new algorithms coming into play, Zhang, Patuwo, and Hu (1998) used NN for forecasting. Recently, artificial neural networks (ANN) have played attention enhancing devotions in the field of time series predicting (Hamzaçebi (2008), Zhang (2003), Kihoro, Otieno, and Wafu (2004)). In particular, ANNs have the advantage of not assuming the statistical distribution followed by the values, being able of proficiently capturing non-linearities. That is, they are self-adaptive (Zhang, Patuwo, and Hu (1998), Zhang (2003)).

A particular class of ANN is the called recurrent neural network (RNN). RNN allows to use previous outputs as inputs, this allows the model to retain information about the past, making it very efficient for time series. This has been shown by Dorffner (1996). In this article, the author highlights the power of RNN for forecasting compared with standard linear models.

For this reason, a lot of authors interested in forecasting economic time series have compared linear models with nonlinear models. In economic time series, in the short run, the time series is expected to behave more or less the same way it has been behaving up to this point, but on the other hand, if we are interested in forecasting at a big window, then is when chaos and instability appear, meaning that nonlinear relations may arise, making it more appealing to use ANN as they are capable of identifying these turning points as they do not assume a linear relationship of inputs and outputs.

This was shown by a recent paper of the Bank of England, Joseph et al. (2021). In this paper they run a horse race for forecasting inflation among different horizons comparing the performance of linear and nonlinear algorithms. The results support our hypothesis that NN and other nonlinear Machine Learning algorithms are useful for forecasting at a longer horizon given that, the longer the horizon, the more likely it is to find the chaos and instability Brock et al. (1991) talked about. And, as previously exposed, these turning points are hard to spot with linear relationships of inputs and outputs, while they might be easier to spot with nonlinear models, like SVM or NN.

One of the problems of RNN is the long-term dependency. To illustrate it with an example, some decisions are made taking into account information that happened way back in the past. RNN struggle to keep this very old information threfore limiting its forecasting power. In order to solve this problem, Hochreiter and Schmidhuber (1997) introduced the LSTM in the paper Long Short-Term Memory. This is the reason why a lot of authors use this type of RNN when forecasting any kind of time series.

### 3 Data

To study the monetary transmission mechanism we used a balanced panel of monthly US data which spans the period of January 1959 through March 2021. In order to explain the monetary transmission mechanism, the VAR literature relies on different indicators such as output and income, industrial new orders and turnover, retail sales and turnover, building permits, employment, consumption, price indices, exchange rates, short- and long-term interest rates, stock price indices, money and credit quantity aggregates, balance of payments and external trade, confidence indicators, and some foreign variables such as output, prices, interest rates, and stock markets from other countries used as proxies for external real, nominal and monetary influences.

On this paper we relied on four of the macroeconomic indicators mentioned above; output, price indices, employment and rates. For each indicator we chose a time series as its proxy. We can see the proxies used for each variable below.

- *output*: Industrial Production (IP)
- *price indices*: Consumer Price Index (CPI)
- *unemployment*: Unemployment rate (UR)
- *rates*: US Federal Fund Rates (FFR)

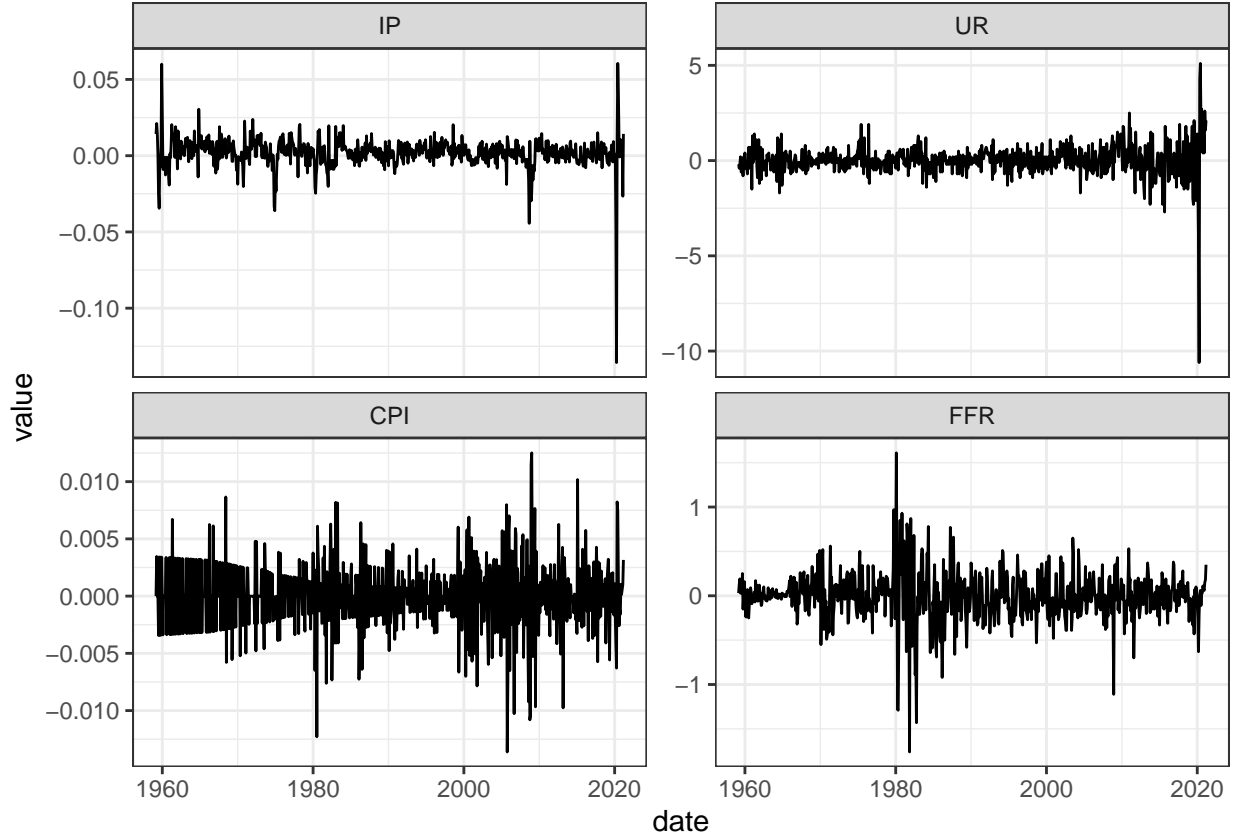
Note that we use IP rather than GDP as a proxy for output because for the case of the GDP only quarterly data was available.

As consequence of working with monthly data, we worked 745 observations, each with the four variables previously mentioned.

The data is taken from the Economic Research department of the Federal Reserve Bank of Saint Louis (FRED). In a recent paper, they have created a monthly database for economic research (FRED-MD), a macroeconomics data base which consists of 134 monthly US indicators. The data is automatically updated each month and can be downloaded freely, facilitating the replication of results and comparing results.

Another strength of using the FRED-MD is the fact that the data is already preprocessed. To see how it has been preprocessed see the Methodology section.

In the graph below we can find the time series already preprocessed by the FRED-MD.



	IP	UR	CPI	FFR
min	-0.14	-10.60	-0.01	-1.76
max	0.06	5.10	0.01	1.61
mean	0.00	0.02	0.00	-0.00
sd	0.01	0.81	0.00	0.27
var	0.00	0.66	0.00	0.08
median	0.00	0.00	-0.00	0.00
IQR	0.01	0.70	0.00	0.28

Table 1: Summary statistics

## 4 Pre Covid vs Post Covid Summary Stats

## 5 Methodology

In conventional Vector Autoregression (VAR) dependencies of any system variable on past realizations of itself and its covariates are modelled through linear equations. This corresponds to a particular case of the broader class of Deep Vector Autoregressions investigated here and will serve as the baseline for our analysis.

	IP	UR	CPI	FFR
min	-0.04	-2.70	-0.01	-1.76
max	0.06	2.50	0.01	1.61
mean	0.00	0.01	0.00	-0.00
sd	0.01	0.63	0.00	0.27
var	0.00	0.40	0.00	0.08
median	0.00	0.00	-0.00	0.00
IQR	0.01	0.63	0.00	0.28

Table 2: Summary statistics pre Covid

	IP	UR	CPI	FFR
min	-0.14	-10.60	-0.01	-0.63
max	0.06	5.10	0.01	0.35
mean	-0.00	0.68	0.00	0.01
sd	0.05	4.04	0.00	0.23
var	0.00	16.31	0.00	0.05
median	0.01	1.60	0.00	0.06
IQR	0.01	2.20	0.00	0.10

Table 3: Summary statistics post Covid

## 5.1 Vector Autoregression

Let  $\mathbf{y}_t$  denote the  $(K \times 1)$  vector of variables at time  $t$ . Then the VAR( $p$ ) with  $p$  lags and a constant deterministic term is simply a linear system of stochastic equations of the following form:

$$\mathbf{y}_t = \mathbf{c} + \mathbf{A}_1 \mathbf{y}_{t-1} + \mathbf{A}_2 \mathbf{y}_{t-2} + \dots + \mathbf{A}_p \mathbf{y}_{t-p} + \mathbf{u}_t \quad (1)$$

The matrices  $\mathbf{A}_m \in \mathbb{R}^{K \times K}$ , where  $m \in \{1, \dots, p\}$ , contain the reduced form coefficients and  $\mathbf{u}_t \in \mathbb{R}^{K \times 1}$  is a vector of errors for which  $\mathbb{E} \mathbf{u}_t$ ,  $\mathbb{E} \mathbf{u}_t \mathbf{u}_t^T = \Sigma$  and  $\mathbb{E} \mathbf{u}_t \mathbf{u}_s^T = \mathbf{0}$  for all  $t \neq s$ . We refer to (1) as the **reduced form** representation of the VAR( $p$ ) because all right-hand side variables are predetermined (Kilian and Lütkepohl 2017).

We can restate (1) more compactly as

$$\mathbf{y}_t = \mathbf{A} \mathbf{Z}_{t-1} + \mathbf{u}_t \quad (2)$$

where  $\mathbf{A} = (\mathbf{c}, \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_p) \in \mathbb{R}^{K \times (Kp+1)}$  and  $\mathbf{Z}_{t-1} = (1, \mathbf{y}_{t-1}^T, \dots, \mathbf{y}_{t-p}^T)^T \in \mathbb{R}^{(Kp+1) \times 1}$ . The expression in (2) demonstrates that the VAR( $p$ ) can be considered as a **seemingly unrelated regression** (SUR) model composed of individual regressions with common regressors (Greene 2012). In fact, it is useful to note for our purposes that the VAR( $p$ ) can be estimated efficiently through equation-by-equation OLS regression. In particular, it follows from (2) that

$$y_{it} = c_i + \sum_{m=1}^p \sum_{j=1}^K a_{jm} y_{jt-m} + u_{it} \quad , \quad \forall i = 1, \dots, K \quad (3)$$

which corresponds to the key modelling assumption that at any point in time  $t$  any time series  $i \in 1, \dots, K$  is just a weighted sum of past realizations of itself and all other variables in the system. This assumption makes the estimation of VAR( $p$ ) processes remarkably simple. Perhaps more importantly, the assumption of linearity also greatly facilitates inference about VARs.

For implementation purposes it is generally more useful to estimate the VAR( $p$ ) through one single OLS regression. To this end let  $\tilde{\mathbf{A}} = \mathbf{A}^{-1}$  and note that (2) can be restated even more compactly as

$$\mathbf{y} = \mathbf{Z}\tilde{\mathbf{A}} + \mathbf{u}_t \quad (4)$$

with  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_T)^T \in \mathbb{R}^{T \times K}$  and  $\mathbf{Z} \in \mathbb{R}^{T \times (Kp+1)}$ . Then the closed form solution for OLS is simply  $\tilde{\mathbf{A}} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$  and hence

$$\mathbf{A} = \mathbf{y}^T \mathbf{Z} (\mathbf{Z} \mathbf{Z}^T)^{-1} \quad (5)$$

## 5.2 Deep Vector Autoregression

We propose the term Deep Vector Autoregression to refer to the broad class of Vector Autoregressive models that use deep learning to model the dependences between system variables through time. In particular, as before we let  $\mathbf{y}_t$  denote the  $(K \times 1)$  vector that describes the state of system at time  $t$ . Consistent with the conventional VAR we assume that each individual time series  $y_{it}$  can be modelled as a function of lagged realizations of all variables  $y_{jt-p}$ ,  $j = 1, \dots, K$ ,  $m = 1, \dots, p$ . More specifically we have

$$y_{it} = f_i(\mathbf{y}_{t-1:t-p}; \theta) + v_{it} \quad , \quad \forall i = 1, \dots, K \quad (6)$$

where  $\mathbf{y}_{t-1:t-p} = \{y_{jt-m}\}_{j=1, \dots, K}^{m=1, \dots, p}$  is the vector of lagged realizations,  $f_i$  is a variable specific mapping from past lags to the present and  $\theta$  is a vector of parameters. While in the conventional VAR above we assumed that the multivariate process can be modelled as a system of linear stochastic equations, our proposed Deep VAR( $p$ ) can similarly be understood as a system of potentially highly non-linear equations. As we argued earlier, Deep Learning has been shown to be remarkably successful at learning mappings of arbitrary functional forms (Goodfellow, Bengio, and Courville 2016).

Note that the input and output dimensions in (6) are exactly the same as in the conventional VAR( $p$ ) model (equation (3)):  $f_i$  maps from  $\mathbf{y}_{t-1:t-p} \in \mathbb{R}^{Kp \times 1}$  to a scalar. Our proposed plain-vanilla approach to Deep VARs diverges as little as possible from the conventional approach: it boils down to simply modelling each of the univariate outcomes in (6) as a deep neural network. We can restate this approach more compactly as

$$\mathbf{y}_t = \mathbf{f}(\mathbf{y}_{t-1:t-p}; \theta) + \mathbf{v}_t \quad (7)$$

where  $\mathbf{f}(\cdot) = (f_1(\cdot), f_2(\cdot), \dots, f_K(\cdot))^T \in \mathbb{R}^{K \times 1}$  is just the stacked vector of mappings to univariate outcomes described in (6).

The notation in (7) gives rise to a more unified and general approach to Deep VARs that would treat the whole process as one single dynamical system to be modelled through one deep neural network  $\mathbf{g}$ :



$$\mathbf{y}_t = \mathbf{g}(\mathbf{y}_{t-1:t-p}; \theta) + \mathbf{v}_t \quad (8)$$

This approach is in fact proposed and investigated by Verstyuk (2020) in his upcoming publication. We decided to go with the approach in (7) for two reasons: firstly, the link to conventional VARs is made abundantly clear through this implementation and, secondly, we found that the equation-by-equation approach produces good modelling outcomes and is relatively easy to implement using state-of-the-art software.

Finally, note that if  $f_i$  in (3) is assumed to be linear and additive for all  $i = 1, \dots, K$  then we are back to the conventional VAR( $p$ ). This illustrates the point we made earlier that the linear VAR( $p$ ) is just a particular case of a Deep VAR( $p$ ). Since the model described in equations (6) and (7) is less restrictive but otherwise consistent with the conventional VAR framework, we expect that it outperforms the traditional approach towards modelling multivariate time series processes.

## 5.3 Deep Neural Networks - a whistle-stop tour

So far we have been speaking about deep learning in rather general terms. For example, above we have referred to our model of choice for learning the mapping  $f_i : \mathbf{y}_{t-1:t-p} \mapsto y_{it}$  as a **deep neural network**. The class of deep neural networks can further be roughly divided into **feedforward neural networks** and **recurrent neural networks**. As the term suggests, the latter is generally used for sequential data and therefore our preferred model of choice. Nonetheless, below we will begin by briefly exploring feedforward neural networks first. This should serve as a good introduction to neural networks more generally and (even though we have not tested this empirically) there is good reason to believe that even Deep VARs using feedforward neural networks perform well.

### 5.3.1 Deep Feedforward Neural Networks

The term **deep feedforward neural network** or **multilayer perceptron** (MLP) is used to describe a broad class of models that are composed of possibly many functions that together make up the directed acyclical graph. The functions  $f_i(\cdot)$  - sometimes referred as layers  $\mathbf{h}_i$  - are chained together hierarchically with the first layer feeding forward its outputs to the second layer and so on (Goodfellow, Bengio, and Courville 2016). Applied to our case, an MLP with  $H$  hidden layers can be loosely defined as follows:

$$f_i(\mathbf{y}_{t-1:t-p}; \theta) = f_i^{(H)} \left( f_i^{(H-1)} \left( \dots f_i^{(1)} (\mathbf{y}_{t-1:t-p}) \right) \right) \quad (9)$$

The depth of the MLP is defined by the number of hidden layers  $H$ , where, generally speaking, deeper networks are more complex.

Need a reference here.

The desired outputs of any  $f_i^{(h)}$  that will serve as inputs for  $f_i^{(h+1)}$  cannot be inferred from the training data  $\mathbf{y}_{t-1:t-p}$  ex-ante, which is where the term **hidden** layer stems from. Each  $f_i^{(h)}$  is typically valued on a vector of hidden units, each of them receiving a vector of inputs from  $f_i^{(h-1)}$  and returning a scalar that is referred to as activation value. This approach is inspired by neuroscience, hence the term **neural** network (Goodfellow, Bengio, and Courville 2016).

### 5.3.2 Deep Recurrent Neural Networks

**Recurrent neural networks** (RNN) are based on the idea of persistent learning: a continuous process that evolves gradually and at each step uses information about its prior states instead of continuously reinventing itself and starting from scratch. To this end, RNNs develop the basic concepts underlying feedforward neural networks by incorporating feedback loops. Formally the loop is typically made explicit as follows

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t; \theta) \quad (10)$$

where  $\mathbf{h}_t \in \mathbb{R}^{N \times 1}$  corresponds to the hidden state of the dynamical system at time  $t$  that the RNN learns (Goodfellow, Bengio, and Courville 2016), and  $N$  corresponds to the number of hidden units in each hidden layer, known as the width of the layer. In the given context we have that  $\mathbf{x}_t = \mathbf{y}_{t-1:t-p}$  as specified in (7). Given some random initial hidden state vector  $\mathbf{h}_0$  the RNN updates parameters sequentially at each time step  $t$  as follows

$$\begin{aligned} \mathbf{a}_t &= \mathbf{b} + \mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{h}_{-1} \\ \mathbf{h}_t &= \tanh(\mathbf{a}_t) \\ \hat{\mathbf{y}}_t &= \mathbf{c} + \mathbf{V}\mathbf{h}_t \end{aligned} \quad (11)$$

where  $\mathbf{b} \in \mathbb{R}^{N \times 1}$  and  $\mathbf{c} \in \mathbb{R}^{K \times 1}$  are vectors of constants (biases),  $\tanh$  is the hyperbolic tangent activation function and  $\mathbf{W}$ ,  $\mathbf{U}$ ,  $\mathbf{V}$  are coefficient matrices, where  $\mathbf{W}, \mathbf{U} \in \mathbb{R}^{N \times N}$  and  $\mathbf{V} \in \mathbb{R}^{K \times N}$ . Note that to simplify the notation we have omitted the layer index in (11): to be specific,  $\mathbf{h}_t$  really represents  $\mathbf{h}_t^{(H)}$  (the ultimate hidden layer),  $\mathbf{h}_{-1}$  stands for  $\mathbf{h}_t^{(H-1)}$  (the penultimate layer). Finally, at each step  $t$  the first layer  $\mathbf{h}_t^{(0)}$  of the forward propagation corresponds to  $\mathbf{y}_{t-1:t-p}$ .

A shortfall of generic recurrent neural networks is that they fail to capture long-term dependencies. More specifically, if parameters are propagated over too many stages in a simple RNN it typically suffers from the problem of **vanishing gradients** (Goodfellow, Bengio, and Courville 2016). Fortunately, there exist effective extensions of the RNN, most notably the long short-term memory (LSTM), which is our model of choice for Deep VARs. The key idea underlying LSTMs is to regulate exactly how much information is propagated from one cell state vector  $\mathbf{C}_{t-1}$  to the next  $\mathbf{C}_t$  through the introduction of so called sigmoid gates:

“The LSTM [has] the ability to remove or add information to the cell state, carefully regulated by structures called gates. Gates are a way to optionally let information through.” — (???)

These regulating gate layers include a **forget gate**  $\mathbf{f}_t$ , an **input gate**  $\mathbf{i}_t$  and a **output gate**  $\mathbf{o}_t$ . Each of them are vector-valued sigmoid functions whose elements  $\mathbf{f}_{it}, \mathbf{i}_{it}, \mathbf{o}_{it}$  are bound between 0 and 1. Their individual purposes are implied by their names: faced with  $\mathbf{h}_{t-1}$  and  $\mathbf{y}_{t-1:t-p}$ , the forget gate regulates how much of each individual unit in  $\mathbf{C}_{t-1}$  is retained. Then the input gate regulates which units of  $\mathbf{C}_{t-1}$  should be updated and to what candidate values  $\tilde{\mathbf{C}}_{t-1}$ . Using the previous two steps the actual update is performed according to the following rule

$$\mathbf{C}_t = \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{C}}_{t-1} \quad (12)$$

where  $\odot$  indicates the element-wise product. Finally, the output gate acts like a filter on  $\mathbf{C}_t$ : the new hidden state is computed as  $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{C}_t)$  where as before we use the hyperbolic tangent as our activation function.<sup>1</sup> Formally, we can summarize the LSTM neural network underlying our Deep VAR framework as follows:

$$\begin{aligned}
\mathbf{f}_t &= \sigma(\mathbf{b}_f + \mathbf{W}_f \mathbf{h}_{t-1} + \mathbf{U}_f \mathbf{h}_{-1}) \\
\mathbf{i}_t &= \sigma(\mathbf{b}_i + \mathbf{W}_i \mathbf{h}_{t-1} + \mathbf{U}_i \mathbf{h}_{-1}) \\
\mathbf{o}_t &= \sigma(\mathbf{b}_o + \mathbf{W}_o \mathbf{h}_{t-1} + \mathbf{U}_o \mathbf{h}_{-1}) \\
\mathbf{C}_t &= \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tanh(\mathbf{b}_C + \mathbf{W}_C \mathbf{h}_{t-1} + \mathbf{U}_C \mathbf{h}_{-1}) \\
\mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{C}_t) \\
\hat{\mathbf{y}}_t &= \mathbf{c} + \mathbf{V} \mathbf{h}_t
\end{aligned} \tag{13}$$

which is best understood when read from top to bottom. Once again we have simplified the notation by omitting the layer index in (11). The same notation as before applies.

## 5.4 Model selection

There are at least two important modelling choices to be made in the context of conventional VARs. The first choice concerns properties of the time series data itself, in particular the order of integration and cointegration. The second choice is about the lag order  $p$ . In order to arrive at appropriate decisions regarding these choices the VAR literature provides a set of guiding principles. We propose to apply these same principles to the Deep VAR, firstly because they are intuitive and simple and secondly because treating both models equally to begin with allows for a better comparison of the two models at the subsequent modelling stages.

### 5.4.1 Stationarity

When working with time series we are generally concerned about stationarity. Broadly speaking stationarity ensures that the future is like the past and hence any predictions we make based on past data adequately describe future outcomes. In order to state stationarity conditions in the VAR context it is convenient to restate the  $K$ -dimensional VAR( $p$ ) process in companion form as

$$\mathbf{Y}_t = \begin{pmatrix} \mathbf{c} \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \mathbf{A} \mathbf{Y}_{t-1} + \begin{pmatrix} \mathbf{u}_t \\ 0 \\ \vdots \\ 0 \end{pmatrix} \tag{14}$$

where  $\mathbf{Y}_t = (\mathbf{y}_t^T, \dots, \mathbf{y}_{t-p+1}^T)^T \in \mathbb{R}^{Kp \times 1}$  and  $\mathbf{A} \in \mathbb{R}^{Kp \times Kp}$  is referred to as the companion matrix (Kilian and Lütkepohl 2017). Stationarity of the VAR( $p$ ) follows from stability: a VAR( $p$ ) is stable if the effects of shocks to the system eventually die out. Stability can be assessed through the system's autoregressive roots or equivalently by looking at the eigenvalues of the companion matrix  $\mathbf{A}$  (Kilian and Lütkepohl 2017). In particular, for the VAR( $p$ ) in (14) to be stable we condition that the  $Kp$  eigenvalues  $\lambda$  that satisfy

---

<sup>1</sup>For a clear and detailed exposition see (???).

$$\det(\mathbf{A} - \lambda \mathbf{I}_{Kp}) = 0$$

are all of absolute value less than one. Stability implies that the first and second moments of the VAR( $p$ ) process are time-invariant, hence ensuring weak stationarity (Kilian and Lütkepohl 2017).

A straight-forward way to deal with stationarity of VARs is to simply ensure that the individual time series entering the system are stationary. This usually involves differencing the time series until they are stationary: for any time series  $y_i$  that is integrated of order  $I(\delta)$ , there exists a  $\delta$ -order difference that is stationary. An immediate drawback of this approach is the loss of information contained in the levels of the time series. Modelling approaches that take into account cointegration of individual time series can ensure system stationarity and still let individually non-stationary time series enter the system in levels (Hamilton 1994).

### 5.4.2 Lag order

The VARs lag order  $p$  can to some extent be thought of as the persistency of the process: past innovations that still affect outcomes in time  $t$  happened at most  $p$  periods ago. From a pure model selection perspective we can also think of additional lags in terms of additional regressors that add to the model's complexity. From that perspective, choosing a lower lag order corresponds to a form of regularization as it pertains to a more parsimonious model.

Various strategies have been proposed to estimate the true or optimal lag order  $p$  empirically (Kilian and Lütkepohl 2017). Among the most common ones are sequential testing procedures and selection based on information criteria. The former involves sequentially adding or removing lags - **bottom-up** and **top-down** testing, respectively - and then testing model outcomes in each iteration. A common point of criticism of sequential procedures is that the order tests matters.

Insert reference to Lütkepohl (2005, Section 4.2.3)

Here we will focus on selection based on information criteria, which to some extent makes the trade-off between bias and variance explicit (Kilian and Lütkepohl 2017). In particular, it generally involves minimizing information criteria of the following form

$$C(m) = \log(\det(\hat{\Sigma}(m))) + \ell(m) \tag{15}$$

where  $\hat{\Sigma}$  is just the sample estimate of the covariance matrix of errors and  $\ell$  is a loss function that penalizes high lag orders. In particular, we have that our best estimate of the optimal lag order  $p$  is simply

$$\hat{p} = \arg \min_{m \in \mathcal{P}} C(m) \tag{16}$$

where  $\mathcal{P} = [m_{\min}, m_{\max}]$ . We will consider all of the most common functional choices for (15).

### 5.4.3 Neural Network Architecture

As mentioned above, a NN is therefore characterized for having a very large number of parameters and therefore making it relatively easy to overfit the data. This is a concern specially when the NN is trained in a relatively small dataset. In time series the amount of data to train the model is quite limited, therefore having so many parameters may result in the NN learning the noise of the function, potentially leading to overfitting which in turn results in fairly poor performance results of the model in the test set. Therefore some regularization must be added into the model in order to avoid overfitting.

As proposed in Srivastava et al. (2014), one way to reduce overfitting can be fitting a lot of NN in the data and the averaging its predictions. This is not feasible due to computational problems. What one can do instead is to randomly drop some units in each layer. This is also known as dropout, this technique simulates the training of multiple NN with different architectures in the same data set. By using dropout, we add some noise into the model. This avoids a certain layer trying to adapt to a mistake made by a previous one, hence potentially leading to complex adaptations which results in overfitting. Dropout therefore can be used as a regularization technique that makes the model more robust.

The way we implement dropout in the Deep VAR model is by specifying the probability to which a node is removed from the layer. This means that the information that arrives to that unit is forgotten.

As another effort to avoid overfitting, we use the same architecture for each layer of the NN. That is, we do not allow each layer to have different number of units, nor different activation functions nor different dropout rates, which could potentially result in the model learning the noise. In turn the Deep VAR model is characterized for having equally sized layers of size  $N$ .

Finally, with respect to how the LSTMs underlying the Deep VAR are compiled, the Adam optimization algorithm Kingma and Ba (2014) is used. This algorithm can be used instead of the more traditional stochastic gradient descent to update network weights. There are several reasons to use this algorithm that are particularly appealing, among them we have its straightforward implementation, that it is computationally efficient and that the hyper-parametrization has an intuitive interpretation and typically requires little tuning. Adam is different to classic stochastic gradient descent given that in the latter learning rates and always the same for all weight updates and hence does not change during the training.

## 6 Empirical results

We now proceed to benchmark the proposed Deep VAR model against the conventional VAR using out macroeconomic time series data. To begin with, we compare both models in terms of their in-sample fit. For this part of the analysis the models will be strictly run under the same framing conditions. Due to the RNN's capacity to essentially model any possible function  $f_i(\cdot)$  the Deep VAR dominates the VAR in this realm. We investigate during what time periods the outperformance of the Deep VAR is particularly striking to gain a better understanding of when and why it pays off to relax the linearity constraint.

These findings with respect to in-sample performance provide some initial evidence in favor of the Deep VAR. But since a reduction in modelling bias is typically associated with an increase in

variance, we are particularly interested in benchmarking the models with respect to their out-of-sample performance. To this end we split our sample into train and test subsamples. We then firstly benchmark the models in terms of their pseudo out-of-sample fit. Finally we also look at model performance with respect to  $n$ -step ahead pseudo out-of-sample forecasts.

The final part of this section relaxes the constraint on the framing conditions. In particular, we investigate how hyperparameter tuning with respect to the neural network architecture and lag length  $p$  can improve the performance of the Deep VAR.

## 6.1 In-sample fit

For this first empirical exercise both models are trained on the full sample. We have decided to include the post-Covid sample period despite the associated structural break, since it serves as interesting point of comparison. The optimal lag order as determined by the Akaike Information Criterion is  $p = 6$ , where we used a maximum possible lag of  $p_{\max} = 12$  corresponding to one year. The LSTMs underlying the Deep VAR model are composed of  $H = 2$  that count  $N = 100$  hidden units each. The dropout rate is set to  $p = 0.5$ .

To assess the fit of our models we use the root mean squared error (RMSE) as our preferred loss function. Figure 1 shows the cumulative RMSE of both the VAR model and Deep VAR model for each of the time series over the whole sample period. The first thing we can observe is that the RMSE of the Deep VAR is consistently flatter than the RMSE of the VAR. With respect to in-sample performance, the Deep VAR the VAR throughout the entire time period of the experimental analysis and for all of the considered variable. This empirical observation seems to confirm our expectation that the vector autoregressive process is characterized by important non-linear dependencies across time and variables that the conventional VAR fails to capture.

Figure 1 is especially useful to asses in which specific periods the Deep VAR model achieves better modelling outcomes than the VAR model. From the very beginning and across variables, we observe that the increase in cumulative loss for the VAR model is greater than for the Deep VAR model. The US economy during 1960s was influence by John F. Kennedy’s introduction of **New Economics**, which was informed by Keynesian ideas and characterized by increasing levels of inflation, a reduction in unemployment and output growth. The change in government certainly corresponded to a regime switch with respect to the economy (???) and in that sense it is interesting to observe that the Deep VAR appears to be doing a better job at capturing the underlying changes. The 1970s can be broadly thought of as a continuation of New Economy and loosely defined as a period of stagflation. The Deep VAR continues to outperform the VAR during that period.

The first truly interesting development we can observe in Figure 1 coincides with the onset of the Volcker disinflation period. Following years of sustained CPI growth, Paul Volcker set the Federal Reserve on course for a series of interest rate hikes as soon as he became chairperson of the central bank in August 1979. The shift in monetary policy triggered fundamental changes to the US economy and in particular the key econmic indicators we are analysing here throughout the 1980s (???). Despite this structural break, the increase in the cumulative RMSE of the Deep VAR remains almost constant during this decade for most variables. The performance of the VAR on the other hand is unsurprisingly poor over the same period, in particularly so for the CPI and the Fed Funds Rate, which arguably were the two variables most directly affected by the change in policy. The Deep VAR also clearly dominates the VAR with respect to the output related variables (IP) and to a lesser extent unemployment. These findings indicate that changes to the

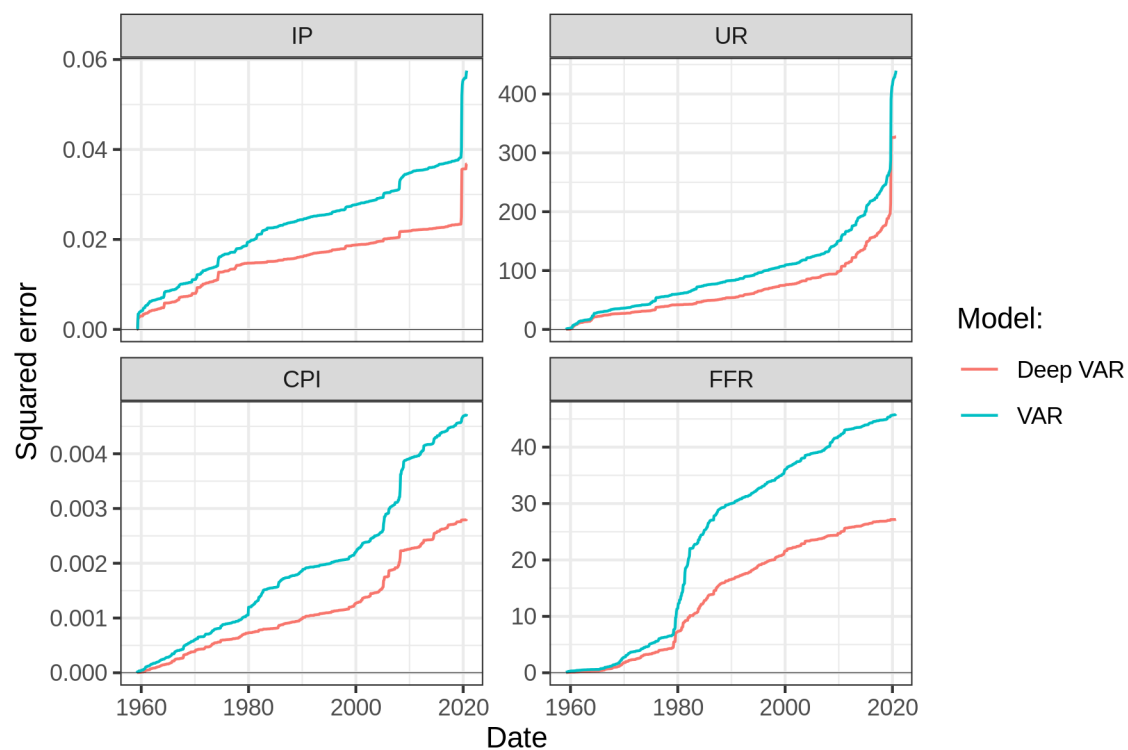


Figure 1: Comparison of cumulative loss over the entire sample period for conventional VAR and proposed Deep VAR.

monetary transmission mechanism in response to sudden policy shifts are not well captured by a linear-additive vector autoregressive model. Instead they appear to unfold in a high-dimensional latent state space, which the Deep VAR by its very construction is designed to learn.

Following the Volcker disinflation period, Figure 1 does not reveal any clear outperformance of either of the models during the 1990s. Interestingly the dot-com bubble has little affect on either of the models, aside from a small pick-up in cumulative loss with respect to the CPI for both models. With all that noted, the Deep VAR still continuously outperforms the VAR since evidently its cumulative loss increases at a lower pace alltogether.

As the Global Financial Crisis unfolds around 2007 the pattern we observed for the Volcker disinflation reemerges, albeit to a lesser extent: there is a marked jump in the difference between the cumulative loss of the VAR and the Deep VAR, in particular so for the CPI, the Fed Funds rate and industrial production. The gap for all these variables continues to widen during the aftermath of the crisis. The Deep VAR once again does a better job at modelling the changes that the dynamical system undergoes: post-crisis US monetary policy was characterized by very low interest rates, low levels of inflation as well as the introduction of a range of non-conventional monetary policy tools including quantitative easing and forward guidance.

Finally, it is also interesting to observe how both models perform in response to the unprecedented exogenous shock that Covid-19 constitutes. Both models incur huge errors with respect to both IP and UR - the two series most significantly affected by Covid. Evidently though, the magnitude of the errors is somewhat larger for the VAR than for the Deep VAR. This, once again seems to confirm our hypothesis that the Deep VAR model captures important non-linear dependences across time and variables that the conventional VAR fails to capture.

As a sanity check we also visually inspected the distributional properties of the model residuals for the full-sample fit. The outcomes are broadly consistent across models: while for some variables residuals are clearly not Gaussian, we see no evidence of serial autocorrelation of residuals. Visualizations can be found in section @ref(#resids) of the appendix.

## 6.2 Out-of-sample fit

In order to assess if the Deep VAR's outperformance is a consequence of overfitting, we now repeat the previous exercise, but this time we train the models on a subsample of our data. The training sample spans from marzo, 1959 to octubre, 2008, whereas the test data goes from noviembre, 2008 to marzo, 2021. This corresponds to training the model on 80 percent of the data and retaining the remaining 20 percent for testing purposes. The optimal lag order for the training subsample is  $p = 7$  where we use the same criterion and maximum lag order as before.

Tables ?? shows the Root Mean Squared Error (RMSE) for the in-sample and the out-of-sample predictions of both the VAR model and the Deep VAR model. We can see that the RMSE for the Deep VAR outperforms the one for the conventional VAR for both the training data and the test data and for all time series. The fifth column of the table shows us the ratio between the RMSEs of the Deep VAR and the VAR: the lower the ratio, the better the Deep VAR compared to the VAR. With respect to the training sample, the RMSE of the Deep VAR model is consistently less than 75% of that of the conventional VAR reflecting to some extent the results of the previous sections. Turning to the test data, there is no evidence that the Deep VAR is more prone to overfitting than the VAR. For both industrial production and unemployment, the Deep VAR yields an RMSE that



is around half the size of that produced by the VAR. For inflation and interest rate predictions the outperformance on the test data is less striking, but still fairly significant.

### 6.3 Forecasts

Up until now we have been assessing the 1-step ahead predictions of both models. In our context these predictions can be thought of as 1-month ahead nowcasts from a practical perspective. Since real-time nowcasts have grown in popularity during recent years, the results so far should be of great interest to central bankers and other practitioners. Nonetheless, there is typically also great interest in time series forecasts at longer horizons. We therefore briefly introduce  $n$ -step ahead pseudo out-of-sample forecasts in this section and revisit them again further below.

Forecasts are produced recursively both for the VAR and the Deep VAR. Specifically, we use the models we trained on the training data to recursively predict one time period ahead, concatenate the predictions to the training data and repeat the process.<sup>2</sup>

In the following table, we look more specifically at this comparison. In particular, we compute the Forecast RMSE for the VAR method and for the Deep-VAR, and we also compute the correlation of the forecast with the actual value of each time series.

Variable	VAR FRMSE	Deep-VAR FRMSE	VAR correlations	Deep-VAR correlations
IP	0.01870	0.01602	-0.30409	-0.65279
UR	0.85984	0.82785	-0.10093	0.27425
CPI	0.00946	0.00708	-0.33567	0.07823
FFR	0.52321	0.39161	-0.55935	0.01161

As we can see in the table, the Forecast RMSE of the Deep-VAR is lower than the one for the VAR for unemployment rate and for the federal funds rate, and it is the same for industrial production and inflation. Hence, the Deep-VAR outperforms the VAR in terms of forecasts. It is also of interest to analyze the correlation between the forecasts and the actual values. In the case of the VAR, there is, for all time series, a negative correlation, that is, when the time series evolves in one direction, the VAR forecast evolves in the opposite direction. However, this is not true for the Deep-VAR. For industrial production, is certainly true that the Deep-VAR forecast has a highly negative correlation with the actual values, but achieves the same forecasting performance than the one for the VAR. However, for the rest of time series, the Deep-VAR has a positive correlation, that is, the forecast generally evolves in the same direction as the actual values. Note that this positive correlation cannot be too strong given that the forecasts at some point in time return to the mean, and hence, have no variability.

### 6.4 Hyper parameter tuning

---

<sup>2</sup>Note that for the Deep VAR an alternative approach would be to work with a different output dimension for the underlying neural networks.

# References

- Bernanke, Ben. 1990. "The Federal Funds Rate and the Channels of Monetary Transnission." National Bureau of Economic Research.
- Brock, William A, David Arthur Hsieh, Blake Dean LeBaron, William E Brock, and others. 1991. *Nonlinear Dynamics, Chaos, and Instability: Statistical Theory and Economic Evidence*. MIT press.
- Cortes, Corinna, and Vladimir Vapnik. 1995. "Support-Vector Networks." *Machine Learning* 20 (3). Springer: 273–97.
- Dorffner, Georg. 1996. "Neural Networks for Time Series Processing." In *Neural Network World*. Citeseer.
- Enders, Walter. 2008. *Applied Econometric Time Series*. John Wiley & Sons.
- Fix, E, and J Hodges. 1951. "An Important Contribution to Nonparametric Discriminant Analysis and Density Estimation." *International Statistical Review* 3 (57): 233–38.
- Friedman, Milton, and Anna Jacobson Schwartz. 2008. *A Monetary History of the United States, 1867-1960*. Princeton University Press.
- Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- Greene, William H. 2012. "Econometric Analysis 7th Ed (International)." Pearson.
- Hamilton, James Douglas. 1994. *Time Series Analysis*. Princeton university press.
- Hamzaçebi, Coşkun. 2008. "Improving Artificial Neural Networks' Performance in Seasonal Time Series Forecasting." *Information Sciences* 178 (23). Elsevier: 4550–9.
- Ho, Tin Kam. 1995. "Random Decision Forests." In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, 1:278–82. IEEE.
- Hochreiter, Sepp, and Jürgen Schmidhuber. 1997. "Long Short-Term Memory." *Neural Computation* 9 (8). MIT Press: 1735–80.
- Joseph, Andreas, Eleni Kalamara, George Kapetanios, and Galina Potjagailo. 2021. "Forecasting Uk Inflation Bottom up."
- Kihoro, J, RO Otieno, and C Wafula. 2004. "Seasonal Time Series Forecasting: A Comparative Study of Arima and Ann Models." *African Journal of Science; Technology (AJST)*.
- Kilian, Lutz, and Helmut Lütkepohl. 2017. *Structural Vector Autoregressive Analysis*. Cambridge University Press.
- Kingma, Diederik P, and Jimmy Ba. 2014. "Adam: A Method for Stochastic Optimization." *arXiv Preprint arXiv:1412.6980*.
- Kydland, Finn E, and Edward C Prescott. 1982. "Time to Build and Aggregate Fluctuations." *Econometrica: Journal of the Econometric Society*. JSTOR, 1345–70.
- Lucas Jr, Robert E. 1976. "Econometric Policy Evaluation: A Critique." In *Carnegie-Rochester Conference Series on Public Policy*, 1:19–46. North-Holland.
- McCulloch, Warren S, and Walter Pitts. 1990. "A Logical Calculus of the Ideas Immanent in Nervous Activity." *Bulletin of Mathematical Biology* 52 (1-2). Springer: 99–115.

- Romer, Christina D, and David H Romer. 1989. “Does Monetary Policy Matter? A New Test in the Spirit of Friedman and Schwartz.” *NBER Macroeconomics Annual* 4. MIT Press: 121–70.
- Sims, Christopher A, and others. 1986. “Are Forecasting Models Usable for Policy Analysis?” *Quarterly Review* 10 (Win). Federal Reserve Bank of Minneapolis: 2–16.
- Srivastava, Nitish, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. “Dropout: A Simple Way to Prevent Neural Networks from Overfitting.” *The Journal of Machine Learning Research* 15 (1). JMLR. org: 1929–58.
- Verstyuk, Sergiy. 2020. “Modeling Multivariate Time Series in Economics: From Auto-Regressions to Recurrent Neural Networks.” *Available at SSRN 3589337*.
- Zhang, G Peter. 2003. “Time Series Forecasting Using a Hybrid Arima and Neural Network Model.” *Neurocomputing* 50. Elsevier: 159–75.
- Zhang, Guoqiang, B Eddy Patuwo, and Michael Y Hu. 1998. “Forecasting with Artificial Neural Networks:: The State of the Art.” *International Journal of Forecasting* 14 (1). Elsevier: 35–62.