

Exploring the Potential of Dynamic Factor Models for High-Dimensional Time Series Forecasting

Bachelor Thesis

January 16, 2023

Author: Linus Christian Wolff

Supervisor: Dr. Toni Stocker

Abstract

This thesis investigates the application of Dynamic Factor Models for forecasting purposes in the context of high-dimensional time series data. After introducing the theoretical framework of the Dynamic Factor Model and its estimation procedure by means of Principal Component Analysis, we simulate a variety of Data Generating Processes that analyze the forecasting performance of a DFM through a Monte Carlo simulation. Our results imply that the forecasting accuracy is highly dependent on the characteristics of the underlying data. Our Dynamic Factor Model approach outperforms the benchmark model substantially only if the variation within the data is dominantly driven by a small number of underlying factors.

Contents

1	Introduction	2
1.1	Literature Review	4
2	Methodology	6
2.1	Dynamic Factor Model Framework	6
2.2	Principal Component Analysis Framework	7
2.3	Estimation of Factor Model by means of PCA	8
2.4	Determining the number of Factors	9
2.5	Forecasting methods	10
3	Monte Carlo Simulation Study	12
3.1	Data Generating Processes	13
3.2	Simulation Results - DGP 1	15
3.3	Simulation Results - DGP 2	19
3.4	Simulation Results - DGP 3	22
3.5	Simulation Results - DGP 4	25
4	Conclusion	28
5	Discussion and possible further research	30
6	References	31

1 Introduction

As information technology improved and continues to improve, the volume and availability of financial and economic data has increased immensely in recent years and decades. That this trend will continue for the foreseeable future is inevitable and will provide both opportunities and dangers to researchers, policymakers and others. Our ability to both collect and process much larger amounts of data than just years ago, theoretically allows us to exploit much more information, hopefully leading to more precise economic forecasting and analysis. Low dimensional settings, where the number of variables is much smaller than the sample size, have historically been the norm in Statistics. Since then, dealing with datasets constantly increasing in size has become common in industry and academia. With the proliferation of “big data”, a lot of the classical methods used by researchers in the context of time series econometrics have run into technical challenges.

While usually low-dimensional Vector Autoregressions show reasonably accurate forecasting performance and had become a standard tool in macroeconometric analysis in the 1990’s, they do confront researchers with two challenges. First, adjusting VAR’s to a high-dimensional context, in which there are possibly more variables than observations, poses a computational bottleneck, as an increase in the number of variables leads to a squared increase in the number of parameters. Second, the non trivial task of deciding which variables to include and exclude in a VAR framework lies with the researcher.

Dynamic Factor Models are a dimension reduction technique that has gained popularity since the early 2000’s because it offers a number of advantages here, especially in a scenario with an abundance of data series. Underlying them is the idea that a large part of the variation contained in economic time series with a large number of variables can be summarised by a much smaller number of unknown factors. Generally speaking, when dealing with factor models, the researcher is neither restricted by the number of variables to include, nor does he have to make any choices as to which variables to include. That being said they usually also lack the interpretability of low dimensional VAR or similar models.

Moreover, the early hope that high-dimensional Dynamic Factor Models would be a breakthrough in forecasting macroeconomic variables seems to have been misplaced. The literature does suggest however that they have resulted in substantial forecasting improvements, not only in simulations but also for measures of real economic activity. Comparing the fore-

casting performance of a Dynamic Factor Model to a classic AR benchmark model will be the main focus of this thesis. Through a simulation study we will try to analyze if, and under which characteristics inherent in the underlying data, Dynamic Factor Models possess improved forecasting performance when compared to our benchmark model.

It should be stated regardless, that a particularly useful application of Dynamic Factor Models has turned out to be the monitoring of economies in real time. The majority of economic indicators are released intermittently, at different frequencies and the number of series has gotten large. Factor Model based techniques, namely the construction of indices and “nowcasting”, have become an important macroeconomic monitoring tool.

Principal Component Analysis is a convenient method for the estimation of the factors in a Dynamic Factor Model that is among the most widely used methods in the literature, especially for high-dimensional data. As mentioned above, standard techniques such as Vector Autoregressions, have had to contend with the “curse of dimensionality”. As the number of variables approaches or even exceeds the number of points in time, linear regression methods cannot be used anymore, due to the lack of degrees of freedom. The Factor Model approach turns this “curse” upside down, because it can handle an arbitrary amount of variables and its precision can even improve as the number of variables increases. This is one of the main reasons why Dynamic Factor Models has received a lot of interest by macroeconomists. The empirical literature focuses on macroeconomic forecasting and the construction of economic indicators most prominently. In an effort to exploit the growing dimensionality of time series data, both public and private institutions now commonly use them for these purposes as well.

In this thesis, Section 1.1 will discuss some of the important literature in both the theoretical and the empirical research on Dynamic Factor Models. Section 2.1 will lay out the methodology of the Factor Model framework, followed by a brief explanation of Principal Component Analysis in Section 2.2. The estimation of a DFM by means of Principal Component Analysis follows in Section 2.3. We will subsequently present a method for determining the number of factors and an explanation of the chosen forecasting procedure. A Monte Carlo experiment and its results, based on multiple Data Generating Processes, is conducted in Section 3. Section 4 comes to a conclusion, while Section 5 outlines possible directions of further research.

1.1 Literature Review

Over the past three decades, Dynamic Factor Models have become one of the leading methods for modeling and forecasting in the context of high-dimensional economic time series.

Early theoretical work on Dynamic Factor Models was done in Sargent and Sims (1977), who proposed the exact factor model. Additionally Sargent and Sims (1977) formulated a single factor model that managed to explain the majority of the variance within multiple monthly economic activity measures. The approximate factor model has its roots in Chamberlain and Rothschild (1983). Coming from the exact factor model, it relaxes the assumption of diagonality of the covariance matrix of the idiosyncratic component, therefore allowing for cross-correlation of the idiosyncratic component.

This formulation of the approximate factor model makes it more suited to economic data and high dimensional data more generally, as the assumptions of the exact factor model are often too strong in these cases. This idea was then further generalized in Forni et. al (2000), which allows for a dynamic representation of the common component and both serially and cross-correlated idiosyncratic components.

Utilizing Dynamic Factor Models for prediction purposes was shown in Stock and Watson (1998, 2002a, 2002b) using large macroeconomic datasets. A comparison of the use of dynamic and static factors in a forecasting context can be found in Boivin and Ng (2005). They find the method using dynamic factors employed in Stock and Watson (2002a) to provide systemically better forecasting performance compared to using static factors and other methods. An alternative approach to forecasting is layed out in Forni et al. (2005), reporting improvements over the results of Stock and Watson (2002b).

Multiple methods for estimating the factors have been discussed in the literature, including Principal Component Analysis, Maximum Likelihood estimation, Kalman filtering and Bayesian methods. For high-dimensional factor models, the majority of the literature focuses on estimation by means of Principal Component Analysis. An alternative approach is taken in Doz et al. (2012), making use of Maximum Likelihood methods. This thesis will stick to estimation using Principal Component Analysis throughout.

Connor and Korajczyk (1986) suggests that when N is larger than T , the factor model can be estimated by applying PCA to the sample covariance of the data. When N and T are large, Bai (2003) shows that the PCA estimator of the common component is asymptotically Gaussian, even when heteroskedasticity and correlation in the idiosyncratic component is present.

The use of Dynamic Factor Models for predicting real and nominal variables using empirical data has been demonstrated in Stock and Watson (2002b) for the US, in Marcellino, Stock and Watson (2003) for the euro area and in Schumacher and Dreger (2004) for Germany.

They have also been utilized for the construction of economic indicators. One of the most prominent examples is the EuroCOIN, an indicator of economic activity in the euro area. First developed by Altissimo et al. (2001), the EuroCOIN is the common component of euro area GDP, estimated through the use of dynamic principal component analysis using high dimensional time series data. A similar indicator, the Chicago Fed National Activity Index (CFNAI) is methodologically based on earlier research in Stock and Watson (1999) about inflation forecasting, and is simply the first static principal component of a panel of 85 macroeconomic time series.

2 Methodology

2.1 Dynamic Factor Model Framework

As mentioned before, Dynamic Factor Models are a dimension reduction technique that aims at identifying the underlying factors describing a large part of the variation within a large number of correlated time series. It is assumed that a small number of unobserved factors explain the correlation of a large number of variables, where $r \ll N$. The common component depends on these unobserved factors, while the idiosyncratic component encapsulates the series-specific movements not captured by the underlying factors, as well as measurement error. The idiosyncratic component is assumed to be uncorrelated with the common component, but can be weakly correlated to the other idiosyncratic components of other series. We assume the data matrix X to be standardized throughout, i. e. it has a mean of 0 and a standard deviation of 1.

Let T be the number of time series observations and N be the number of observed variables. We define an approximate factor model for $i = 1, \dots, N$ and $t = 1, \dots, T$ as

$$\begin{aligned} x_{it} &= \lambda_{i1}F_{1t} + \dots + \lambda_{ir}F_{rt} + \xi_{it} \\ &= \lambda_i'F_t + \xi_{it} \\ &= C_{it} + \xi_{it}. \end{aligned} \tag{2.1}$$

We define λ_i as the factor loadings for variable i corresponding to the r common factors F_t , where r is the true number of factors underlying the model. The common component is represented by $C_{it} = \lambda_i'F_t$, while the idiosyncratic component is ξ_{it} .

Now, let $X_t = (x_{1t}, x_{2t}, \dots, x_{Nt})'$, $F = (F_1, \dots, F_T)'$ and $\Lambda = (\lambda_1, \dots, \lambda_N)'$. Equivalently, in vector form we get

$$X_t = \Lambda F_t + \xi_t \tag{2.2}$$

Lastly, let $\mathbf{X} = (X_1', \dots, X_N')$ be a matrix of observations with dimension $T \times N$. The matrix

representation of the factor model is therefore

$$\mathbf{X} = \mathbf{F}\Lambda' + \xi, \quad (2.3)$$

where \mathbf{F} is $T \times r$, Λ' is $r \times N$ and $\xi = (\xi'_1, \xi'_2, \dots, \xi'_N)$ is $T \times N$.

We allow for a dynamic relationship of the r factors F , by specifying that they follow a Vector Autoregressive process (VAR), which is the second part of the DFM and will be important for forecasting using the DFM. The VAR takes the form of

$$\begin{aligned} F_{1t} &= A_{11}F_{1t-1} + \dots + A_{1r}F_{rt-1} + \eta_{1t} \\ &\vdots \\ F_{rt} &= A_{r1}F_{1t-1} + \dots + A_{rr}F_{rt-1} + \eta_{rt} \end{aligned} \quad (2.4)$$

It is a VAR(1) above, there are generally no restrictions on the number of lags included however. The η 's represent error terms with $E(\eta_{rt}) = 0$, whereas the A 's are unknown VAR coefficients to be estimated.

2.2 Principal Component Analysis Framework

The insight that a variable that is highly correlated with other regressors could be dropped without losing a lot of information, forms the basis for Principal Component Analysis. Particularly when many regressors are present, PCA can help to reduce the number of predictors with a minimal loss of information. By definition, PCA is a dimension reduction technique aiming at explaining the majority of information within a set of correlated variables by its Principal Components.

We will now consider a data matrix X of dimension $n \times p$ which is assumed to have a zero mean, meaning $E(X_j) = 0$ for $j = 1, \dots, p$. The $p \times p$ sample covariance matrix of this data is denoted by Σ_x . The eigenvalues corresponding to Σ_x are denoted as $\hat{\lambda}_1, \dots, \hat{\lambda}_p$ ordered such that $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$. The eigenvectors are $\hat{e}_1, \dots, \hat{e}_p$ corresponding to the eigenvalue with the same subscript, i.e. such that the eigenvector \hat{e}_1 is the one that corresponds to the largest eigenvalue.

$$\begin{aligned}
\hat{PC}_j &= \hat{e}_j^T X \\
\hat{PC}_{ij} &= \hat{e}_j^T X_i = \hat{e}_{j1}x_{i1} + \dots + \hat{e}_{jp}x_{ip}
\end{aligned} \tag{2.5}$$

gives the j -th principal component for the i -th observation. As the above equation shows, the principal components are linear combinations of our initial p variables X_1, \dots, X_p with the following properties:

- The squared weights, the elements of \hat{e}_j , equal 1, such that $\hat{e}_j^T \hat{e}_j = 1$
- The j -th principal component is uncorrelated to the first $j - 1$ principal components, as $Cov(\hat{PC}_j, \hat{PC}_k) = 0$ for $j \neq k$.
- The j -th principal component maximizes the variance of its linear combination, because it holds that $Var(\hat{PC}_j) = \lambda_j$.

To sum up, one first takes the covariance matrix of the data in X . Second, its eigenvalues and eigenvectors get computed. Form a matrix of the normalized eigenvectors columnwise, ordered by the eigenvalues, such that $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$. Multiplying the resulting matrix with X yields the matrix of Principal Components.

2.3 Estimation of Factor Model by means of PCA

Now that we have described the general frameworks of both Dynamic Factor Models and Principal Component Analysis, the estimation procedure of Dynamic Factor Models by means of PCA will be layed out. The approach taken here will follow Stock and Watson (2002a).

As explained in the PCA section, we will be start off by computing the sample covariance matrix (assuming our data has a mean of zero) $\hat{\Gamma}^x = \frac{1}{T} \sum x_t x_t' = \frac{X'X}{T}$ of dimension $N \times N$. The $r \times r$ diagonal matrix \hat{M}^x contains the r largest eigenvalues of $\hat{\Gamma}^x$ in descending order. The corresponding normalized eigenvectors are the columns of the $N \times r$ matrix \hat{V}^x . Estimators for the factor loadings Λ and the factors F are obtained by solving the minimization problem

$$\min_{\Lambda F} \frac{1}{NT} \sum_{i=1}^n \sum_{t=1}^T (x_{it} - \lambda_i' F_t)^2 = \min_{\Lambda F} \frac{1}{NT} \text{tr} \left\{ (X - F\Lambda')' (X - F\Lambda') \right\} \tag{2.6}$$

The above objective function minimizes the sum of the squared residuals between the observed data and the estimated factors, which is our standard least squares problem.

There are multiple equivalent ways of solving Equation (2.6), namely in Bai (2003), Forni et al. (2009) and Stock and Watson (2002a). We will follow the latter approach. We solve Equation (2.6) by first substituting $F = X\Lambda(\Lambda'\Lambda)^{-1} = \frac{X\Lambda}{N}$ and solving for Λ subject to $\frac{\Lambda'\Lambda}{N} = I_r$ before finally solving for F by linear projection. After substituting F in Equation X, we have

$$\min_{\Lambda} \frac{1}{NT} \text{tr} \left\{ X(I_n - \frac{\Lambda\Lambda'}{N})X' \right\} \quad (2.7)$$

or equivalently

$$\max_{\Lambda} \frac{1}{NT} \text{tr} \left\{ X'X \frac{\Lambda\Lambda'}{N} \right\} = \max_{\Lambda} \frac{1}{N} \text{tr} \left\{ \frac{\Lambda'}{\sqrt{N}} \frac{X'X}{T} \frac{\Lambda}{\sqrt{N}} \right\} \quad (2.8)$$

The assumption $\frac{\Lambda\Lambda'}{N} = I_r$ not only ensures that the columns of Λ are uncorrelated (orthogonal), but also that they are of unit length, or in other words normalized. The estimator $\hat{\Lambda}$ that maximizes above equation must be set so that $\frac{\hat{\Lambda}}{\sqrt{N}}$ is equal to the matrix of normalized eigenvectors of $\frac{X'X}{NT}$, that corresponds to its r largest eigenvalues.

With \hat{M}^x containing the r largest eigenvalues, it follows that

$$\frac{\hat{\Lambda}'}{\sqrt{N}} \frac{X'X}{NT} \frac{\hat{\Lambda}}{\sqrt{N}} = \frac{\hat{M}^x}{n}, \quad (2.9)$$

which can be rewritten as

$$\hat{V}^{x'} \frac{X'X}{NT} \hat{V}^x = \frac{\hat{M}^x}{n}. \quad (2.10)$$

Therefore, it can be shown that $\hat{\Lambda} = \hat{V}^x \sqrt{N}$ is the resulting $N \times r$ estimated loadings matrix, while $\hat{F} = \frac{X\hat{\Lambda}}{\sqrt{N}}$ gives the corresponding $T \times r$ matrix of estimated factors.

2.4 Determining the number of Factors

Our aim is to construct a factor model that describes the comovements within X through a limited number of factors so that $r \ll N$, meaning that r is substantially smaller than N .

With regards to PCA, there are several common approaches for determining the appropriate number of principal components. Namely, two methods most commonly used are visual evaluation of scree plots, which shows the ordered eigenvalues of the covariance matrix $\hat{\Sigma}_X$ divided by the sum of the eigenvalues, and evaluation of the cumulative variance explained by the first k principal components directly.

The econometric literature commonly follows the approach taken by Bai and Ng (2002), in which multiple Information Criteria for determining the number of factors r are proposed. We will be using the most widely adopted in the literature, which is their IC_{p2} criterion:

$$IC_{p2}(r) = \frac{1}{NT} \sum_{t=1}^T (X - \Lambda F_t)' (X_t - \Lambda F_t) + r \left(\frac{N+T}{NT} \right) \ln(\min(N, T)). \quad (2.11)$$

This criterion takes the natural log of the least squares objective function mentioned earlier and adds a penalty increasing proportionally with the number of factors. As for other Information criteria, the value of r that minimizes $IC_{p2}(r)$ will be chosen. For computational purposes, we restrict the maximum number of factors to be $r = 8$.

2.5 Forecasting methods

The main forecasting method that we will utilize here is a Factor Augmented Auto Regression. Starting off, the Factor Model in Equation (2.1) will be expanded upon such that

$$x_{it} = \lambda_{i0} + \lambda_{i1}F_{1t} + \dots + \lambda_{ir}F_{rt} + \beta_1 x_{it-1} + \dots + \beta_p x_{it-p} + u_{it}. \quad (2.12)$$

Our forecasting equation assumes that lagged values of x_{it} might be useful predictors besides the factors we have estimated previously. This equation could also be expanded upon further by allowing the lagged values of other series regarded as useful predictors to enter the equation, which would give a Factor Augmented Autoregressive Distributed Lag Model. When forecasting with this approach, a problem arises. Let h be the forecasting horizon. The current values of the factors are unknown for any x_{it+h} where $h > 0$, thus they cannot be used as predictors.

To remedy this, we will be incorporating the iterated forecasting approach in the following way. The factor VAR estimates from Equation (2.4) will replace the current factor values in Equation (2.12), such that a one-step ahead forecast $h = 1$ for x_{it} is

$$\hat{x}_{iT+1} = \hat{\lambda}_{i0} + \hat{\lambda}_{i1}\hat{F}_{1T+1} + \cdots + \hat{\lambda}_{ir}\hat{F}_{rT+1} + \hat{\beta}_1x_{iT} + \cdots + \hat{\beta}_px_{iT-p+1}. \quad (2.13)$$

Here, $\hat{F}_{1T+1}, \dots, \hat{F}_{rT+1}$ are the $h = 1$ forecasts from the factor VAR equation described in Equation (2.4). The estimates of the $\hat{\lambda}$ and $\hat{\beta}$ coefficients are computed using Equation (2.12), where lagged values of x_{it} and $\hat{F}_{1t}, \dots, \hat{F}_{rt}$ are used as regressors. Although we will be focusing on one-step ahead forecasts here, as the name iterated implies, the forecasts for $h > 1$ horizons use the iterated VAR forecasts for the factors and the preceding forecast of x_i .

We will be comparing this approach to a baseline model, which is the univariate autoregressive model. Forecasts based on this AR(p) model take the form

$$\hat{x}_{iT+1} = \hat{\beta}_0 + \hat{\beta}_1x_{iT} + \cdots + \hat{\beta}_px_{iT-p+1}, \quad (2.14)$$

where the lag length p is chosen based on the Akaike information criterion (AIC). The AIC gets computed as

$$AIC(p) = \ln \left[\frac{SSR(p)}{T} \right] + (p+1)\frac{2}{T}, \quad (2.15)$$

where $SSR(p)$ is the sum of squared residuals of the estimated AR(p). The chosen model is the one that minimizes the AIC.

3 Monte Carlo Simulation Study

The goal of this thesis is to evaluate the forecasting performance of Dynamic Factor Models incorporated into a Factor Augmented Autoregression as shown above. Multiple Data Generating Processes will be setup and simulated. Subsequently, out of sample forecasting will be performed, utilizing both of the aforementioned forecasting equations. Namely, the FAAR approach as developed in Stock and Watson (2002a) will be compared to our simple AR(p) benchmark forecast.

Two measures will be taken to compare the forecasts. First, we make use of the Mean Squared Forecast Error of the out of sample forecasts that will be computed, comparing the MSFE's of the two models we are evaluating. Second, the Diebold-Mariano Test will be performed. The Diebold-Mariano test is a statistical test used to compare the accuracy of two or more forecasting models. It compares the out-of-sample forecast errors of different models by constructing a test statistic based on the mean squared forecast error differences between the models.

The MSFE for an h step ahead forecast of X_i , the $MSFE$ will be computed as

$$MSFE_{X_i} = \frac{1}{T-S} \sum_{t=S}^T (X_{i,t+h} - \hat{X}_{i,t+h})^2, \quad (3.1)$$

where S indicates the last observation included in the data used for estimation, while $\frac{1}{T-S}$ indicates the number of observations used for pseudo out of sample forecasts.

In order to check whether the difference between the MSFE's of our two forecasting models are statistically significant, we will make use of the Diebold-Mariano Test. The test statistic of the Diebold-Mariano Test is given by

$$DM = \frac{\bar{d}}{\sqrt{\frac{\hat{V}(d)}{T}}} \stackrel{approx}{\sim} N(0, 1) \quad (3.2)$$

where T is the length of the forecast window. The sample mean of the loss differential of the squared forecast errors of two series $d = (\hat{X}_{it} - X_{it})^2 - (\hat{X}_{kt} - X_{kt})^2$ is denoted by \bar{d} , while $\hat{V}(d)$ is the sample variance of d over the forecasting horizon $t = 1, \dots, T$, such that

$$\hat{V}(d) = \frac{1}{T-1} \sum_{t=1}^T (d_t - \bar{d})^2. \quad (3.3)$$

The DM test statistic is then compared to a critical value from a standard normal distribution to determine whether there is a significant difference in forecast accuracy between the models. If the test statistic exceeds the critical value, we reject the following null hypothesis:

H0: *The forecast errors of the DFM model do not differ significantly from the forecast errors of the benchmark model.*

More specifically, as the Diebold-Mariano test follows the $N(0, 1)$ distribution asymptotically and we will be using the two-sided version of it, the null hypothesis will be rejected for values outside $[-1.96, 1.96]$ at the standard $\alpha = 0.05$ level, concluding that there is a statistically significant difference in forecast accuracy between the two models.

The following simulations aim at gaining some insight into the relative forecasting performances of the previously mentioned forecasting techniques. First of all, we will be evaluating the general (relative) forecasting accuracies of the proposed models. In addition, the impact of altering a number of parameters incrementally will be assessed.

3.1 Data Generating Processes

All results computed and presented in this section are based on implementations in R.

For $i = 1, \dots, N$, $t = 1, \dots, T$ and $j = 1, \dots, r$, we generate the data in this simulation setup according to the following data generating process:

$$\mathbf{X} = \mathbf{F}\Lambda' + \xi, \quad (3.4)$$

$$F_t = \alpha F_{t-1} + u_t, \quad (3.5)$$

$$\xi_t = \rho \xi_{t-1} + v_t, \quad (3.6)$$

$$\lambda_{ij} = \eta_{ij}. \quad (3.7)$$

The innovations u_{jt} , v_{it} , as well as η_{ij} are i.i.d $N(0, 1)$ random variables. The factors F_t evolve according to a VAR(1) process with a common AR parameter α . Similarly, the

idiosyncratic errors ξ_t are serially correlated through the AR(1) coefficient ρ . The initial values of both the factors and the idiosyncratic errors are drawn from their static distribution.

We consider the following specifications for the free parameters T, N, r, α and ρ :

- The sample size is given by $T = 100$.
- The dimension is given by $N \in \{50, 100, 200\}$.
- The number of true underlying factors is given by $r \in \{2, 4, 6\}$.
- The serial correlation of the factors is given by $\alpha \in \{0, 0.5, 0.8\}$.
- The serial correlation of the idiosyncratic component is given by $\rho \in \{0, 0.2, 0.5\}$.

We split these specifications in the following way. First, an exact factor model is assumed, where we fix ρ at 0, while allowing serial correlation in the factors. The exact factor model assumption of i.i.d idiosyncratic components is then relaxed in the second setup to allow for serial correlation of the idiosyncratic components, resulting in the approximate factor model assumptions.

The last two setups relax the assumptions even further, allowing for both serial- and cross-correlation among the factors, and assuming the idiosyncratic components to either be i.i.d or also possess serial- and cross-correlation. This results in the following four data generating processes displayed in Table 3.1. The specifics of the cross correlation in DGP 3 and DGP 4 will be made precise in later sections.

Table 3.1: The four different data generating processes with different structures of the factors and idiosyncratic errors.

	N	r	alpha	rho
dgp1	{50,100,200}	{2,4,6}	{0,0.5,0.8}	0
dgp2	{50,100,200}	{2,4,6}	{0,0.5,0.8}	{0.2,0.5}
dgp3	100	{2,4,6}	{0,0.5,0.8}	0
dgp4	100	{2,4,6}	0.5	0.2

Each setup is run 100 times. We evaluate three performance criteria for every iteration.

1. For every DGP, the MSFE of the one step ahead forecast errors of the last 20 observations will be computed for the first 10 variables. In other words, if $T = 100$, we compute the (pseudo) out of sample forecast errors for $t = 81, \dots, 100$ for the variables X_i , with $i = 1, \dots, 10$. In this case, the MSFE is as follows:

$$MSFE_X = \frac{1}{10} \sum_{i=1}^{10} \left[\frac{1}{10} \sum_{t=80}^{T=100} (X_{i,t+1} - \hat{X}_{i,t+1})^2 \right]. \quad (3.8)$$

2. As a second measure, the Diebold-Mariano Test will be performed on the forecasts of every variable X_i , with $i = 1, \dots, 10$. We use the resulting p-values as another measure in assessing the forecasting performances of the evaluated models. Specifically, the ratio of p-values that are smaller than 0.05 will be looked at, providing an estimate for the significance in the difference in forecasting accuracy between the evaluated models.
3. Lastly, whether the Information Criterion IC_{p2} proposed above actually determines the number of factors accurately, will be measured by the relative frequency that the number of chosen factors k equals the true number of factors r .

The following models are considered for prediction:

DFM: A dynamic factor model with a VAR for the factors that uses a FAAR (Factor Augmented Auto Regression) forecasting equation. The lag-length of the AR component is set to $p = 4$. The lag-length for the VAR is chosen by AIC and the number of factors is chosen by the information criterion of Section 2.4.

AR: A univariate AR(p) model with the lag-length chosen by AIC.

3.2 Simulation Results - DGP 1

We start off by keeping ρ fixed at $\rho = 0$, meaning the idiosyncratic errors are cross sectionally and serially i.i.d. This is the assumption made in the exact factor model, which is arguably not satisfied when dealing with macroeconomic data. However it does mean that the explanatory power of the factors is large in relation to the idiosyncratic errors, as they possess no correlation in any direction. Furthermore we only allow the factors to be serially correlated through the common AR coefficient α , not cross correlated however.

In the following, the MSFE results will be presented in relation to the MSFE of the AR model, meaning values smaller than 1 indicate better performance and values larger than 1 indicate worse performance in relation to the benchmark AR model.

The overall performance of the two models in this simulation is summarized in Table 3.2. The MSFE of the Dynamic Factor Model is smaller than its benchmark AR model counterpart. We conclude that it does offer some improvement over the benchmark forecasts, while acknowledging that the improvement is minimal.

Table 3.2: relative MSFE in relation to the benchmark
AR model

	DFM	AR
MSFE	0.96	1

In order to judge the relative forecasting performance more precisely, we will take a look at how the specific parameters impact the models forecasting accuracies.

Table 3.3 portrays how the relative MSFE behaves dependent on adjusting the amount of serial correlation present. It is evident that an increase in the serial correlation α entering into the factors leads to an increase in relative forecasting accuracy compared to the benchmark. This is what would be expected, as an $\alpha = 0$ implies no correlation among the factors, essentially letting each factor be an i.i.d white noise process. In other words, at $\alpha = 0$ there is no underlying factor structure at all, resulting in no performance gain against the benchmark. This behavior is therefore expected and encouraging, because it indicates a correct model specification, even if the gain in forecasting performance over the AR model is small.

Table 3.3: relative performance of the models for different
amounts of serial correlation among the factors.

alpha	MSFE	ratio p-values	ratio IC correct
0.8	0.9088	0.1097	0.9378
0.5	0.9689	0.0682	0.9911
0.0	1.0141	0.0629	0.9978

Next, the relative accuracy of the DFM improves slightly for smaller numbers of true underlying factors. This effect is small and will be compared against the other data generating processes.

Table 3.4: relative performance of the models for different numbers of true factors underlying the data.

r	MSFE	ratio p-values	ratio IC correct
2	0.9404	0.0936	0.9989
4	0.9590	0.0792	0.9656
6	0.9923	0.0680	0.9622

In this simulation, the dimension of X does not seem to play a relevant role in our models relative performances, which is in line with the findings of Stock and Watson (2002a). One interesting thing to note however, is that the frequency with which the true number of factors is chosen by the Information Criterion increases slightly as N increases. This is in line with Connor and Korajczyk (1993), which notes that the consistency of the factor estimation improves as $N \rightarrow \infty$.

Table 3.5: relative performance of the models for different dimensions of the data.

N	MSFE	ratio p-values	ratio IC correct
50	0.9698	0.0767	0.9344
100	0.9605	0.0828	0.9922
200	0.9614	0.0813	1.0000

The corresponding Box plots visualize the previously mentioned effects that the free parameters have on the predictive accuracy of the two models. Clearly, the factor structure depending on α possesses the biggest effect.

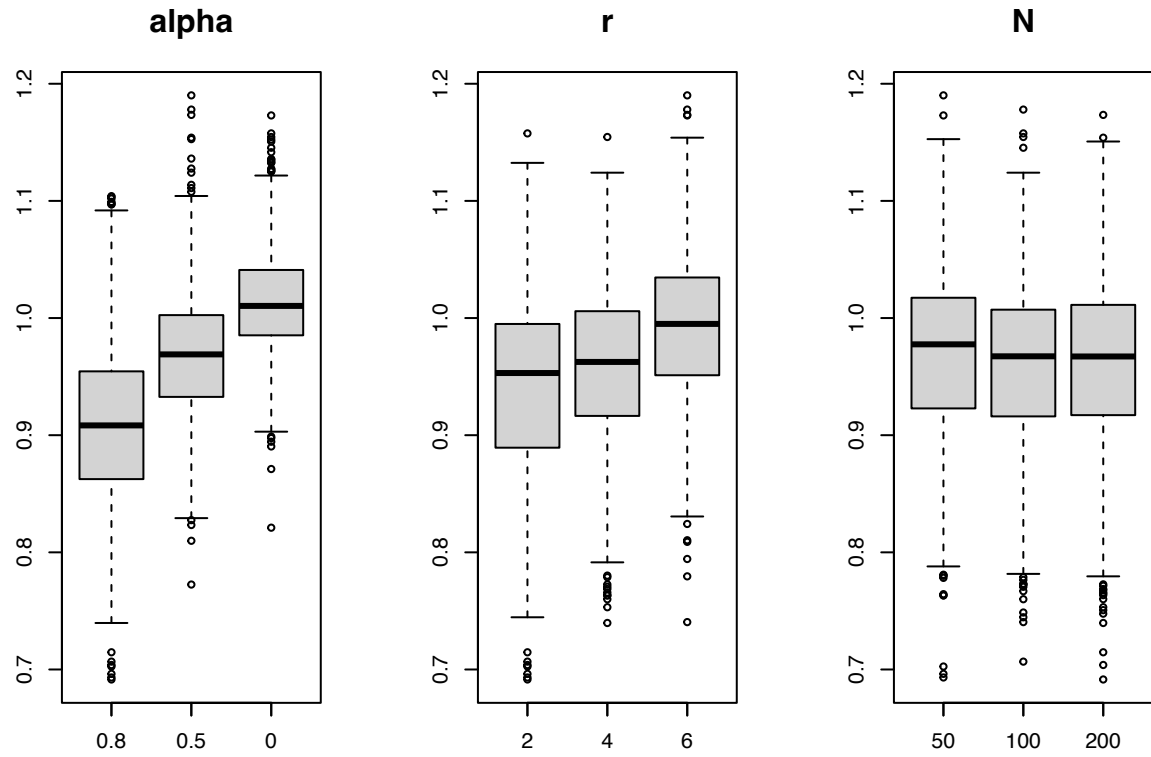


Figure 3.1: RMSE in relation to the benchmark model for every iteration and different parameters α , r and N .

3.3 Simulation Results - DGP 2

We now let $\rho \in \{0.2, 0.5\}$, meaning the idiosyncratic errors are now serially correlated variables. Under this assumption, the factors are less “strong” because a larger part of the covariance of X stems from the structure of the idiosyncratic components.

This is a situation that is often encountered in empirical data. As suggested in Onatski (2012), in which the principal component estimator is found to be inconsistent when the explanatory power of the factors is weak, this should lead to less precise factor estimates and therefore less precise forecasts of the DFM model, when compared to the benchmark.

The relative forecasting performance is displayed in Table 3.6. Indeed, the DFM model even performs slightly worse than the benchmark model in this particular setup.

Table 3.6: relative MSFE in relation to the benchmark AR model.

	DFM	AR
MSFE	1.02	1

The following tables once again disaggregate the impacts that our free parameters have on the models performances in this simulation.

Table 3.7: relative performance of the models for different amounts of serial correlation among the factors.

alpha	MSFE	ratio p-values	ratio IC correct
0.8	1.0087	0.0826	0.9111
0.5	1.0383	0.0792	0.9811
0.0	1.0259	0.0687	0.9961

Table 3.8: relative performance of the models for different numbers of true factors underlying the data.

r	MSFE	ratio p-values	ratio IC correct
2	1.0301	0.0853	0.9956

r	MSFE	ratio p-values	ratio IC correct
4	1.0149	0.0738	0.9567
6	1.0280	0.0714	0.9361

Table 3.9: relative performance of the models for different dimensions of the data.

N	MSFE	ratio p-values	ratio IC correct
50	1.0260	0.0758	0.9017
100	1.0230	0.0777	0.9883
200	1.0238	0.0770	0.9983

Now, the impact of differing the parameter values of α , r and N is insignificant. Two results are in line with the first simulation. First, the ratio of the Information Criterion from Bai and Ng (2002) choosing the true number of factors increases as N increases and decreases as r increases. Second, we conclude from the first two data generating processes that varying $N \in \{50, 100, 200\}$, the dimension of X , essentially had no effect on relative forecasting accuracy. The following two data generating processes will therefore forego varying N , instead fixing it at $N = 100$.

Table 3.10: relative performance of the models for different amounts of serial correlation among the idiosyncratic errors.

rho	MSFE	ratio p-values	ratio IC correct
0.5	1.0609	0.0841	0.9515
0.2	0.9877	0.0696	0.9741

The relative explanatory power of the factors should decrease, resulting in worse performance of the DFM model, as the amount of serial correlation ρ in the idiosyncratic errors increases. The result in Table 3.10 validates this, as it indicates the relative accuracy of the DFM model decreasing for an increasing ρ .

Examining the corresponding Box plots in Figure 3.2 confirms that the two models perform very similarly overall. The only visible trend is the relative accuracy of the DFM model worsening as the factors relatively get weaker.

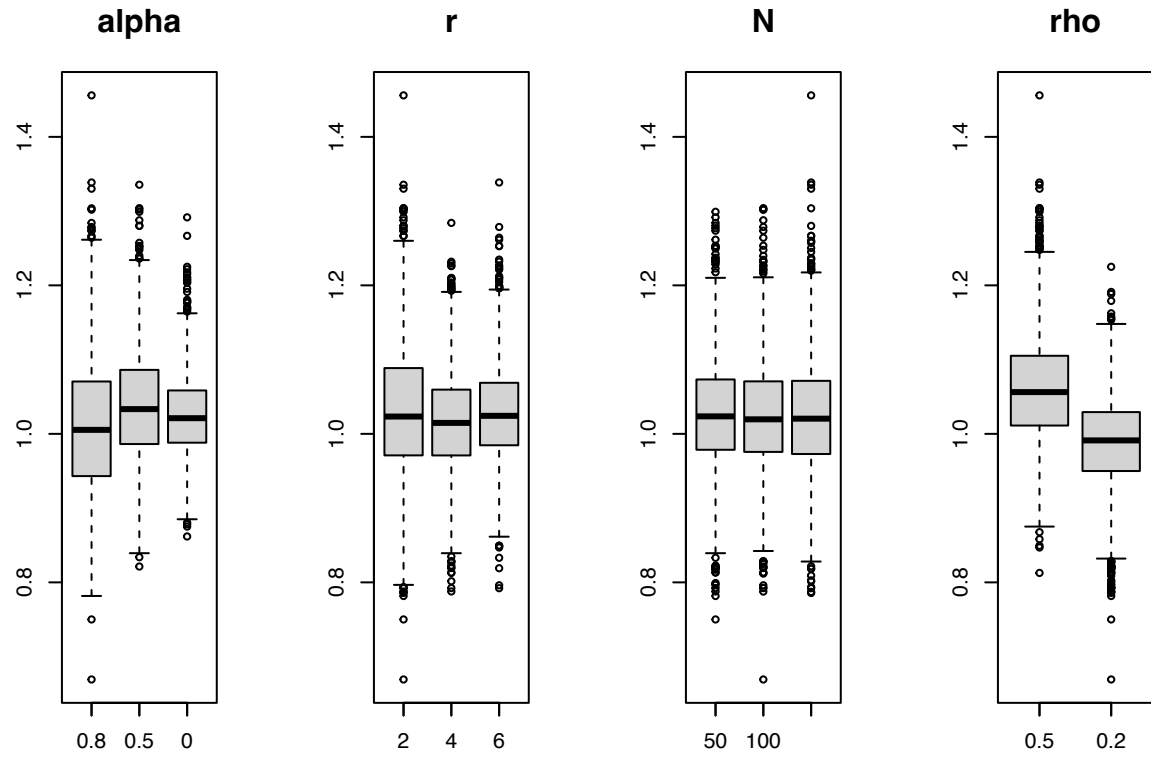


Figure 3.2: RMSE in relation to the benchmark model for every iteration and different parameters alpha, r and N.

3.4 Simulation Results - DGP 3

So far, we have only considered serial correlation among both the factors and the idiosyncratic components. The third data generating process is the extension of the first DGP by now allowing the factors to follow a true VAR(1) process. Therefore, data generating process 3 introduces cross correlation for the factors, while setting $\rho = 0$. That is, as in DGP 1, we assume the idiosyncratic errors to be an i.i.d white noise process.

This setup arguably introduces an empirically plausible structure for the factors, as both dynamics in the time domain, as well as cross correlation between the factors is present. In addition, it theoretically makes full use of the properties of our forecasting equation, the Factor Augmented Auto Regression. It should become clear whether the inclusion of the unobserved factors that we estimate provides a gain in forecasting performance when compared to the simple AR model, especially since the true factors now follow the VAR process that we assume them to follow in the DFM.

The VAR coefficient matrix will be setup as follows, here displayed in the case that $r = 4$:

$$a = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix}$$

Note that the diagonal entries are set as $a_{ii} = \alpha$, the serial correlation coefficient we've established earlier, while all other entries are $a_{ij} \sim U(-0.4, 0.4)$ (uniform distribution) for $i \neq j$.

This might lead to a high correlation between any of the factors, which poses the problem of (perfect) multicollinearity. In that case, our estimations by OLS do not have unique solutions, which would render this setup ill defined. To avoid multicollinearity, we restrict the correlation between any of the factors by imposing that $|corr(F_i, F_j)| < 0.5$ for $i, j = 1, \dots, r$ and $i \neq j$.

As Table 3.11 shows, the Dynamic Factor Model performs substantially better when compared to the benchmark. This result is expected, because this setup gives the true factors the largest amount of relative explanatory power in relation to the idiosyncratic components. Not only does the majority of the exploitable variance underlying this data generating process originate from the factor structure, the DFM and its forecasting equation are also well suited to extract that variance, as we assume underlying factors and that these factors follow

a VAR process. It is therefore unsurprising that this setup yields better relative results for the DFM.

Table 3.11: relative MSFE in relation to the benchmark
AR model

	DFM	AR
MSFE	0.83	1

As Table 3.12 shows, as in the DGP 1, larger amounts of serial correlation of the factors, leads to relative performance gains for the DFM, which leads to the conclusion that both correlation across time, as well as correlation across the variables, gives rise to superior forecasting accuracy for the DFM. On top of that it is noteworthy that the DFM substantially outperforms even in the case of no serial correlation, i. e. when $\alpha = 0$.

Table 3.12: relative performance of the models for different amounts of serial correlation among the factors.

alpha	MSFE	ratio p-values	ratio IC correct
0.8	0.7232	0.2627	0.85
0.5	0.8576	0.1463	1.00
0.0	0.9063	0.1043	1.00

Interestingly, the effect of varying the number of factors r is now contrary to the results of DGP 1. Now, the relative MSFE changes in favor of the DFM when increasing the number of factors. The most obvious reason for this observation would be that the principal components estimates of the factors in combination with the VAR structure of the forecasting equation are better able to extract the variance within the factors as their number grows. The implemented information criterion slightly loses accuracy as the number of factors grows however.

Table 3.13: relative performance of the models for different numbers of true factors underlying the data.

r	MSFE	ratio p-values	ratio IC correct
2	0.9116	0.1097	1.00
4	0.8280	0.1720	0.92
6	0.7475	0.2317	0.93

As the Box plots in Figure 3.3 indicate, the Dynamic Factor Model specification is an improvement over the benchmark AR model in almost all iterations of this setup.

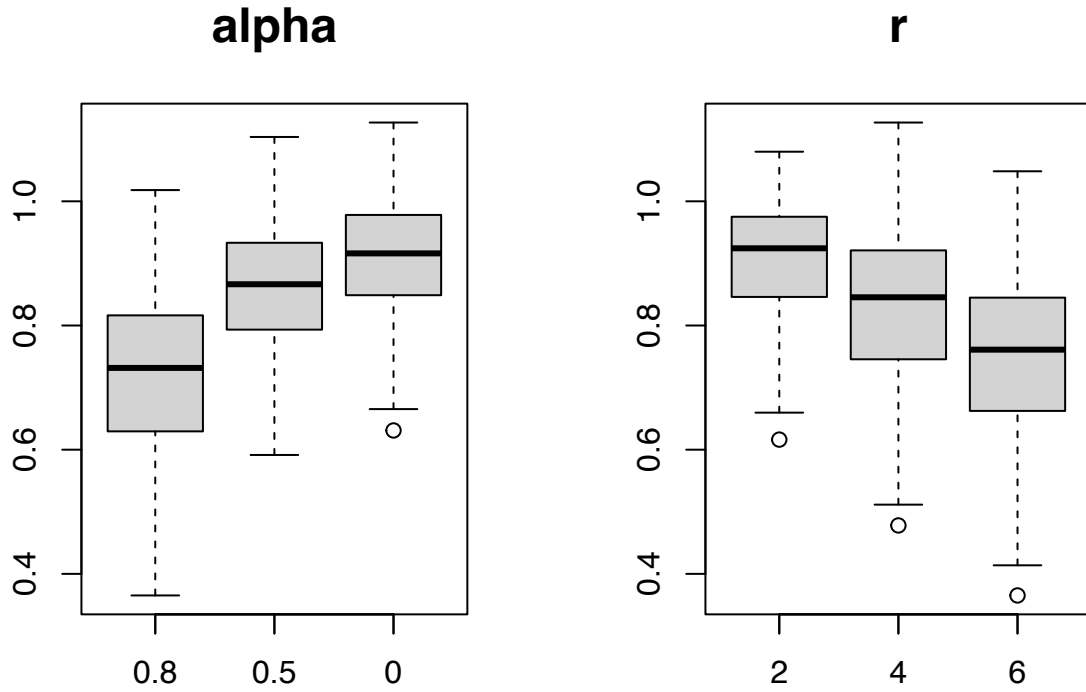


Figure 3.3: RMSE in relation to the benchmark model for every iteration and different parameters α , r and N .

3.5 Simulation Results - DGP 4

Data generating process 4 extends the previous one by now assuming the idiosyncratic components to follow a VAR process as well. Inspired by Krampe and Margaritella (2021), we assume the idiosyncratic components to follow a sparse VAR process, defined such that the VAR coefficient matrix takes the form

$$c = \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1N} \\ c_{21} & c_{22} & \cdots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{N1} & c_{N2} & \cdots & c_{NN} \end{bmatrix},$$

where the diagonal entries are set as $c_{ii} = \rho$, the serial correlation coefficient we've established earlier, while all other entries c_{ij} for $i \neq j$ are drawn from a $U(-0.2, 0.2)$ distribution with probability π and equal to 0 with probability $1 - \pi$. Sparsity is introduced by setting $\pi = 0.2$. For the same reason as before, we restrict the correlation between any of the idiosyncratic components by imposing that $|corr(\xi_i, \xi_j)| < 0.5$ for $i, j = 1, \dots, N$ and $i \neq j$. For computational simplicity, we only consider $\alpha = 0.5$ in this setup.

This data generating process is meant to assess how the performance of the previous approach, which did not allow any correlation among the idiosyncratic components, is impacted as the factor structure explains relatively less variation within the data than in DGP 3.

As would be expected, Table 3.14 indicates that this DGP does indeed perform worse in relative MSFE terms, when compared to the previous specification. It is still outperforming the benchmark model however.

Table 3.14: relative MSFE in relation to the benchmark AR model.

	DFM	AR
MSFE	0.9	1

Once again, the previously noted relationship with regards to the number of factors holds. As r decreases, the performance of the DFM deteriorates. As Table 3.15 indicates, DGP 4 performs worse than DGP 3 for all values of r . The difference in performance decreases as the number of factors increases. This leads to the conclusion that the cross-correlated

nature of the factors increasingly dominates the cross-correlated nature of the idiosyncratic components, measured as their relative explanatory power, as r rises.

Table 3.15: Comparison of the relative performance of the DGP's 3 and 4 for different numbers of true factors underlying the data.

	r	MSFE	ratio p-values	ratio IC correct
dgp3	2	0.9256	0.099	1.00
	4	0.8572	0.155	1.00
	6	0.7900	0.185	1.00
dgp4	2	0.9903	0.090	1.00
	4	0.9120	0.111	1.00
	6	0.8082	0.165	0.98

The relative explanatory power of the factors should decrease, resulting in worse performance of the DFM model, as the amount of serial correlation ρ in the idiosyncratic errors increases. The result in Table 3.10 validates this, as it indicates the relative accuracy of the DFM model decreasing for an increasing ρ .

The corresponding Box plots in Figure 3.4 confirm that the introduction of serial and cross-correlation among the idiosyncratic components worsens the DFM's performance. This is in line with the literature in the sense that the accuracy of DFM based predictions suffers as the importance of the idiosyncratic components grows.

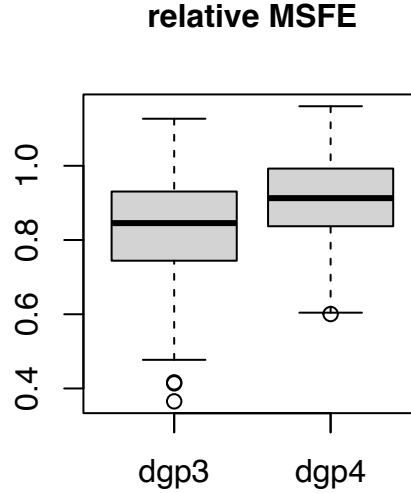


Figure 3.4: RMSE in relation to the benchmark model for data generating processes 3 and 4.

Throughout this simulation study, we have mainly focused on comparing the results of our various data generating process specifications by means of evaluating their respective MSFE as defined in Section 3. We rarely mentioned the ratio with which the Information Criterion from Section 2.4 correctly identified the number of true factors. Suffice it to say that it correctly identified the number of factors in the vast majority of iterations and for every data generating process, without giving us much more information, as there are no clear patterns showing in which cases its performance clearly deteriorates.

Comparing the forecasts based on the Diebold-Mariano Test proved relatively fruitless. When comparing DGP 2 and 3, where the former does not outperform the benchmark at all, while the latter consistently outperforms it, the differences in the ratio that the Diebold-Mariano Test deems the forecasts of the DFM compared to the benchmark to be significantly different is small. This leads to two possible conclusions. Either the performance gains over the benchmark model are not substantial in any of the data generating processes, which is not what the relative MSFE values indicate, or the use of the Diebold-Mariano Test is not expedient in this case. The use and abuse of the test is discussed in Diebold (2015), which mentions that it is not suited for comparing different models in an out of sample forecasting environment, indicating that the implementation of it in this thesis might be one of such “abusive” cases.

4 Conclusion

This thesis gives an overview of the Dynamic Factor Model, its relevance in the econometric literature and its estimation by means of Principal Component Analysis. Furthermore, we cover some of the most widely used techniques for determining the number of factors to incorporate into a DFM, as well as outlining an approach to forecasting within the DFM framework. Namely, a Factor Augmented Autoregression framework closely following the approach taken in Stock and Watson (2002a).

We then set up a Monte Carlo simulation study with the purpose of evaluating the forecasting performance of the Dynamic Factor Model established throughout this thesis in relation to a standard Autoregressive benchmark model. In order to compare the two, we perform an out of sample forecast experiment and evaluate the forecast errors of both models against one another. This is done for a multitude of simulation setups, resulting in four general data generating processes.

Multiple configurations of the data generating processes are considered. We vary the dimension of our generated data and differ the number of true underlying factors. Most importantly, multiple assumptions for the generation of the factors and idiosyncratic components are made. For the factors, depending on the configuration, we evaluate the effects of having no factor structure at all, allowing for serial correlation of the factors, as well as allowing the factors to be serially and cross-correlated. Similarly, the importance of the idiosyncratic component characteristics are assessed. Depending on the configuration we either set them to be independent and identically distributed random variables, serially correlated or following a sparse VAR process.

Overall, our results demonstrate that the Dynamic Factor Model established in this thesis outperforms the standard Autoregressive benchmark model in terms of forecasting performance. Of particular importance is the underlying factor structure we assume. We found that the DFM is able to produce more accurate and efficient forecasts relative to the benchmark model, when the underlying factors were allowed to have serial and cross-correlation. On the other end of the factor structure spectrum, where we assume them to simply be i.i.d random variables, the Dynamic Factor Model predictably performs no better than the benchmark.

Additionally, we find that the DFM frameworks suitability is dependent on the share of the

variation within the data that originates from the factors and the idiosyncratic components respectively. As the factors get weaker, meaning they explain a smaller part of the variance within the data, the relative accuracy of the Dynamic Factor model deteriorates.

We generally find that the forecasting framework of the DFM is sensitive to the underlying structure of the data. It provides a clear improvement in performance compared to the benchmark AR model if the primary force driving the exploitable variation within the data originates in the characteristics of the unobserved factors. In the cases that no unobserved factors are present or the variation in the originating from the idiosyncratic components is dominant, the DFM cannot achieve consistent improvements in performance.

This finding has implications for the application of Dynamic Factor Model approaches for forecasting in an empirical context. The discoveries of this thesis suggest that their use should be limited to situations where it can be strongly argued that a small number of unobserved factors are responsible for most of the variation within a set of variables. Macroeconomics is one of the fields that arguably deals with data that can be characterized this way, explaining why the macroeconometric literature is so vast.

5 Discussion and possible further research

One possible direction for further research is to apply the DFM framework to empirical data such as a high dimensional macroeconomic dataset, as in Bai and Ng (2008) among many others. In the context of both empirical and simulated data, expanding upon the method of Dynamic Factor Models and comparing its performance to other state of the art forecasting methods, would be of great interest. An interesting approach is extending the DFM by combining it with a sparse VAR structure for the idiosyncratic components. Estimating the relationship of the idiosyncratic components and including it in the forecasting regimen is a promising path for performance gains, as Krampe and Margaritella (2021) shows.

More generally, how much data are really needed is a practical question that gets little recognition in this context. Most implementations of the DFM extract the factors from the same high dimensional data. It is not inconceivable however, that different series within a dataset are highly predictable bases on different subsets of the data. Combining “dense” techniques like PCA for factor estimation and “sparse” techniques like the LASSO, elastic net or others for e. g. subset selection, is another interesting direction for further research.

Lastly, it would also be interesting to compare Dynamic Factor Models to suitable machine learning techniques such as Long short-term memory artificial neural networks (LSTM). Comparing the performance of Dynamic Factor Models against these methods could provide insights into the strengths and limitations of both approaches in high dimensional time series forecasting.

6 References

- Altissimo, Filippo, Antonio Bassanetti, Riccardo Cristadoro, Mario Forni, Marc Hallin, Marco Lippi, Lucrezia Reichlin, and Giovanni Veronese. 2001. “EuroCOIN: A Real Time Coincident Indicator of the Euro Area Business Cycle.” *Available at SSRN 296860*.
- Bai, Jushan. 2003. “Inferential Theory for Factor Models of Large Dimensions.” *Econometrica* 71 (1): 135–71.
- Bai, Jushan, and Serena Ng. 2002. “Determining the Number of Factors in Approximate Factor Models.” *Econometrica* 70 (1): 191–221.
- . 2006. “Confidence Intervals for Diffusion Index Forecasts and Inference for Factor-Augmented Regressions.” *Econometrica* 74 (4): 1133–50.
- . 2008. “Forecasting Economic Time Series Using Targeted Predictors.” *Journal of Econometrics* 146 (2): 304–17.
- Chamberlain, Gary. 1983. “Funds, Factors, and Diversification in Arbitrage Pricing Models.” *Econometrica: Journal of the Econometric Society*, 1305–23.
- Connor, Gregory, and Robert Korajczyk. 1993. “A Test for the Number of Factors in an Approximate Factor Model.” *Journal of Finance* 48 (4): 1263–91.
- Connor, Gregory, and Robert A Korajczyk. 1986. “Performance Measurement with the Arbitrage Pricing Theory: A New Framework for Analysis.” *Journal of Financial Economics* 15 (3): 373–94.
- Diebold, Francis X. 2015. “Comparing Predictive Accuracy, Twenty Years Later: A Personal Perspective on the Use and Abuse of Diebold–Mariano Tests.” *Journal of Business & Economic Statistics* 33 (1): 1–1.
- Diebold, Francis X, and Robert S Mariano. 2002. “Comparing Predictive Accuracy.” *Journal of Business & Economic Statistics* 20 (1): 134–44.

- Doz, Catherine, Domenico Giannone, and Lucrezia Reichlin. 2012. “A Quasi-Maximum Likelihood Approach for Large, Approximate Dynamic Factor Models.” *Review of Economics and Statistics* 94 (4): 1014–24.
- Favero, Carlo A., Massimiliano Marcellino, and Francesca Neglia. 2005. “Principal Components at Work: The Empirical Analysis of Monetary Policy with Large Data Sets.” *Journal of Applied Econometrics* 20 (5): 603–20. <https://doi.org/10.1002/jae.815>.
- Forni, Mario, Domenico Giannone, Marco Lippi, and Lucrezia Reichlin. 2009. “Opening the Black Box: Structural Factor Models with Large Cross Sections.” *Econometric Theory* 25 (5): 1319–47.
- Forni, Mario, Marc Hallin, Marco Lippi, and Lucrezia Reichlin. 2000. “The Generalized Dynamic-Factor Model: Identification and Estimation.” *Review of Economics and Statistics* 82 (4): 540–54.
- . 2005. “The Generalized Dynamic Factor Model: One-Sided Estimation and Forecasting.” *Journal of the American Statistical Association* 100 (471): 830–40.
- Krampe, Jonas, and Luca Margaritella. 2021. “Factor Models with Sparse VAR Idiosyncratic Components.” arXiv. <https://doi.org/10.48550/ARXIV.2112.07149>.
- Marcellino, Massimiliano, James H Stock, and Mark W Watson. 2003. “Macroeconomic Forecasting in the Euro Area: Country Specific Versus Area-Wide Information.” *European Economic Review* 47 (1): 1–18.
- Onatski, Alexei. 2012. “Asymptotics of the Principal Components Estimator of Large Factor Models with Weakly Influential Factors.” *Journal of Econometrics* 168 (2): 244–58.
- Sargent, Thomas J, Christopher A Sims, et al. 1977. “Business Cycle Modeling Without Pretending to Have Too Much a Priori Economic Theory.” *New Methods in Business Cycle Research* 1: 145–68.
- Schumacher, Christian, and Christian Dreger. 2004. “Estimating Large-Scale Factor Models for Economic Activity in Germany: Do They Outperform Simpler Models?” *Jahrbücher für Nationalökonomie Und Statistik* 224 (6): 731–50.

- Stock, James H., and Mark W. Watson. 2005. “Implications of Dynamic Factor Models for VAR Analysis.” National Bureau of Economic Research Cambridge, Mass., USA.
- Stock, James H, and Mark W Watson. 1998. “Business Cycle Fluctuations in u.s. Macroeconomic Time Series.” Working Paper 6528. Working Paper Series. National Bureau of Economic Research. <https://doi.org/10.3386/w6528>.
- . 1999. “Forecasting Inflation.” *Journal of Monetary Economics* 44 (2): 293–335.
- . 2002a. “Forecasting Using Principal Components from a Large Number of Predictors.” *Journal of the American Statistical Association* 97 (460): 1167–79.
- . 2002b. “Macroeconomic Forecasting Using Diffusion Indexes.” *Journal of Business & Economic Statistics* 20 (2): 147–62.