

Maximum likelihood estimation of large factor model on datasets with missing data: forecasting euro area GDP with mixed frequency and short-history indicators.

Marta Bańbura¹ and Michele Modugno²

Abstract

This paper deals with the estimation of a large dynamic factor model on datasets with a general pattern of missing data. The framework allows to handle efficiently and in a fairly automatic manner sets of indicators characterized by different publication delays, frequencies and sample lengths. This can be relevant e.g. for young economies for which many indicators are compiled only since recently. We also show how the unexpected part of new data releases is related to the forecast revision. This can be used e.g. to trace the sources of the latter back to the individual series in the case of simultaneous releases. The methodology is applied to the short-term forecasting and backdating of the euro area GDP based on a large panel of monthly and quarterly data.

Keywords: Factor Models, Forecasting, Large Cross-Sections, Missing data, EM algorithm.

JEL classification: C53, E37.

The authors would like to thank Domenico Giannone, Michele Lenza and Lucrezia Reichlin for useful discussions.

The opinions expressed in this paper are those of the authors and do not necessarily reflect the views of the European Central Bank.

¹European Central Bank and ECARES, Université Libre de Bruxelles; email: marta.banbura@ecb.int; contact author;

²European Central Bank and ECARES, Université Libre de Bruxelles; email: michele.modugno@ecb.int.

1 Introduction

This paper deals with the maximum likelihood estimation of a dynamic factor model for large cross-sections with a general pattern of missing data. The framework can be used for extracting information from large datasets including e.g. mixed frequency or short history indicators or for real-time applications in which data arrive with delays and in a nonsynchronous manner. We also show how the unexpected part of a data release is related to the forecast revision and illustrate how this link can be used to understand the sources of the latter in the case of simultaneous releases.

Using a large number of indicators in order to assess the current and the future state of the economy is becoming a standard practice. In this context, factor models have emerged as an effective tool for extracting information from many indicators in a parsimonious way.¹ The fundamental assumption underlying those models is that most of the co-movement of the variables in a given dataset can be summarized by only few factors. This assumption seems to be justified in the case of macroeconomic and financial data. Forecasting macroeconomic variables is one of the domains where factor models show a great potential.²

This paper builds on the approximate dynamic factor model of Giannone, Reichlin, and Small (2008). It has been implemented in several countries and proved to perform well in short-term forecasting, see e.g. Bańbura and Rünstler (2007) for the euro area GDP, Aastveit and Trovik (2007) for Norwegian GDP, Matheson (2007) for New Zealand GDP and inflation. It allows for explicit modelling of the dynamics of the factors³ and for missing data at the end of the sample. This feature can be exploited e.g. to extract information from timely indicators in real-time forecasting.⁴ The framework also allows for a certain degree of cross and serial correlation in the idiosyncratic components.

In Giannone, Reichlin, and Small (2008) the estimation is performed in two steps on a balanced panel.⁵ To establish the relation between the GDP and the factors, the latter are extracted through the principal components analysis on monthly indicators and converted to the quarterly frequency. Therefore the methodology cannot be easily applied to general mixed frequency panels including series of different lengths. The maximum likelihood approach that we propose is able to handle large datasets with an arbitrary pattern of data availability efficiently and in a fairly automatic manner. This is particularly relevant for the euro area or other young economies for which many series have been compiled only since recently (e.g. euro

¹See Forni, Hallin, Lippi, and Reichlin (2000) and Stock and Watson (2002a) for the theoretical foundations.

²See e.g. Bernanke and Boivin (2003); Boivin and Ng (2005); D'Agostino and Giannone (2006); Forni, Hallin, Lippi, and Reichlin (2005, 2003); Giannone, Reichlin, and Sala (2004); Marcellino, Stock, and Watson (2003); Stock and Watson (2002a,b).

³shown to be important in the simulation study of Doz, Giannone, and Reichlin (2005);

⁴in contrast to models based on balanced datasets as in e.g. Forni, Hallin, Lippi, and Reichlin (2000); Stock and Watson (2002a);

⁵The consistency has been proved by Doz, Giannone, and Reichlin (2005).

area Purchasing Managers' (PM) Surveys). Moreover as the series measured at a lower frequency can be interpreted as "high frequency" indicators with missing data, mixed frequency datasets can be easily handled. This can be important for two reasons: first, the information in the indicators sampled at a lower frequency (e.g. consumption, employment) can be used to extract the factors; second, the forecasts of the former can be easily obtained.

Apart from data availability considerations the maximum likelihood estimation is more appealing than the two-step method because it is more efficient and, more importantly, it provides framework for imposing restrictions on the parameters. For example Reis and Watson (2007) impose restrictions on the parameters of the dynamic factor model in order to estimate the pure inflation. Maximum likelihood approach has been used in the classical factor analysis for small datasets (see e.g. Geweke, 1977; Sargent and Sims, 1977; Watson and Engle, 1983). Doz, Giannone, and Reichlin (2006) show that maximum likelihood is consistent, robust and computationally feasible also in the case of large cross-sections. To maximise the likelihood over the high-dimensional parameter space they propose to use the Expectation-Maximisation (EM) algorithm. The EM algorithm was first applied for a dynamic factor model by Watson and Engle (1983) on a small cross-section. They cast the model in a state space form and derive the EM steps in the case without missing data. Shumway and Stoffer (1982) show how to implement the EM algorithm for a state space form with missing data, however only in the case in which the matrix linking the states and the observables is known. We extend those results to the case in which all the parameters of the state space form need to be estimated.

We apply the methodology to short-term forecasting of euro area GDP on the basis of large panel of monthly and quarterly indicators. GDP is an important measure of economic performance, however it is released only with about 2 months delay. In the meantime more timely indicators can provide information on the current state of the economy. For example, Giannone, Reichlin, and Small (2008) claim that the Federal Reserve Bank of Philadelphia Business Outlook Survey is crucial for nowcasting the U.S. GDP due to its early release (cf. Bańbura and Rünstler, 2007, for the case of euro area). We want to examine the effect of quarterly variables and short history monthly series of PM Surveys on the forecast. For example, the PM Surveys are considered to be important soft indicators of the real-activity in the euro area (e.g. Purchasing Managers' Index is analysed in the ECB Monthly Bulletin on a regular basis), however their short sample length is prohibitive for many models that cannot flexibly deal with missing data.

We start by comparing the maximum likelihood approach with the two-step procedure used by Giannone, Reichlin, and Small (2008) and Bańbura and Rünstler (2007) on the benchmark dataset containing 70 monthly "long" indicators dating back to at least to 1993, similar to the one used by e.g Bańbura and Rünstler (2007) or ECB (2008). In the following exercise we augment the benchmark dataset by quarterly and/or short monthly indicators and estimate

the model by maximum likelihood. We find that the results based on different methods and different datasets are comparable. This means that, on one hand, adding quarterly variables and PM Surveys does not lead to forecast accuracy improvements, but on the other hand, our methodology allows us to analyse datasets including the mixed frequency and short history indicators and obtain their forecast in a unified framework. In the following exercise we illustrate how the “news” in the consecutive releases of the industrial production, PM and European Commission Surveys revise the GDP forecast. Finally, we show that the framework can be used for backdating. In particular, the back estimates of GDP are fairly close to the true values.

Related literature includes Camacho and Perez-Quiros (2008) who estimate a small dynamic factor model by maximum likelihood to obtain the real-time estimates of the GDP from monthly indicators. Breitung and Schumacher (2008) forecast GDP from large number of monthly indicators using the EM approach, however they do not exploit the dynamics of the factors. Proietti (2008) estimates a factor model⁶ from a balanced panel by EM for interpolation of expenditure and output components of the GDP and shows how to incorporate relevant accounting and temporary constraints.

As for the interpolation and backcasting, Angelini, Henry, and Marcellino (2006) propose methodology based on large cross-sections. In contrast to their procedure our method exploits the dynamics of the data and is based on maximum likelihood which allows for imposing restrictions and is more efficient for smaller cross-sections.

The paper is organized as follows. Section 2 presents the model, discusses the estimation and explains how the news content can be extracted. Section 3 describes the empirical application. Section 4 concludes. The technical details, data description and robustness checks are provided in the Appendix B.

2 The model

Let $y_t = (y_{1,t}, y_{2,t}, \dots, y_{n,t})$ denote a stationary n -dimensional vector process standardised to mean 0 and unit variance. We assume that y_t follows an approximate factor model given by:

$$y_t = \Lambda f_t + \epsilon_t \quad \epsilon_t \sim N(0, \Psi), \quad (1)$$

where f_t is a $r \times 1$, $r \ll n$ vector of common factors and ϵ_t is the idiosyncratic error, uncorrelated with f_t at all leads and lags. However, contrary to the case of the exact dynamic factor model, the idiosyncratic error can be weakly serially- and cross-correlated (cf. Doz, Giannone, and

⁶a la Watson and Engle (1983) and a la Stock and Watson (2002b);

Reichlin, 2006, for the technical details). Moreover we assume that the common factors f_t follow a stationary VAR process of order p :

$$f_t = A_1 f_{t-1} + A_2 f_{t-2} + \dots + A_p f_{t-p} + u_t \quad u_t \sim N(0, Q). \quad (2)$$

2.1 Estimation

As pointed out in Doz, Giannone, and Reichlin (2006) it is not obvious how to model parametrically the weak cross-correlation of the idiosyncratic component in (1). However they show that the model can be estimated by quasi maximum likelihood, where the miss-specified model is the exact factor model

$$y_t = \Lambda f_t + \xi_t \quad \xi_t \sim N(0, R), \quad (3)$$

with R diagonal and ξ_t serially uncorrelated. They prove that the maximum likelihood estimates of the model given by (1) - (2) and of the miss-specified given by (3) - (2) are asymptotically equivalent. Therefore, in what follows we will derive the estimates under the miss-specified likelihood.

In the main text we set for simplicity $p = 1$, the case of $p > 1$ is discussed in the Appendix. Accordingly, the system given by (3) and (2) can be written in a state space form with the latent factors as states:

$$\begin{aligned} y_t &= \Lambda f_t + \xi_t & \xi_t &\sim N(0, R), \\ f_t &= A_1 f_{t-1} + u_t & u_t &\sim N(0, Q). \end{aligned} \quad (4)$$

A direct maximisation of the likelihood for (4) is computationally not feasible for large n . However, as argued in Doz, Giannone, and Reichlin (2006), the computational complexity can be circumvented by means of the Expectation-Maximisation (EM) algorithm. It offers a solution to problems for which incomplete or latent data yield the likelihood intractable. The essential idea is to write the likelihood as if the data was complete and to “fabricate” the missing data in the expectation step.⁷

Let us collect the parameters of the model in $\theta = (\Lambda, A_1, R, Q)$ and let $l(Y, F; \theta)$ denote the joint log-likelihood of y_t and f_t , $t = 1, \dots, T$. Given the available data Ω_T , the EM algorithm converges towards the maximum likelihood estimates in a sequence of two alternating steps:

1. E-step - the expectation of the log-likelihood conditional on the data is calculated using the estimates from the previous iteration $\theta(r)$:

$$L(\theta, \theta(r)) = E_{\theta(r)} [l(Y, F; \theta) | \Omega_T];$$

⁷The EM algorithm was proposed by Dempster, Laird, and Rubin (1977). For overview see e.g. McLachlan and Krishnan (1996).

2. M-step - the parameters are re-estimated through the maximisation of the expected log-likelihood with respect to θ :

$$\theta(r+1) = \arg \max_{\theta} L(\theta, \theta(r)). \quad (5)$$

From the maximisation of (5) it follows that the estimates in the $r+1$ iteration are given by:

$$A_1(r+1) = \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t f'_{t-1} | \Omega_T] \right) \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_{t-1} f'_{t-1} | \Omega_T] \right)^{-1}, \quad (6)$$

$$\Lambda(r+1) = \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [y_t f'_t | \Omega_T] \right) \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t f'_t | \Omega_T] \right)^{-1}, \quad (7)$$

$$Q(r+1) = \frac{1}{T} \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t f'_t | \Omega_T] - A_1(r+1) \sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_{t-1} f'_t | \Omega_T] \right) \quad (8)$$

and

$$R(r+1) = \text{diag} \left(\frac{1}{T} \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [y_t y'_t | \Omega_T] - \Lambda(r+1) \sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t y'_t | \Omega_T] \right) \right). \quad (9)$$

Watson and Engle (1983) derive the estimates when y_t does not contain missing data. In that case we have

$$\mathbb{E}_{\theta(r)} [y_t y'_t | \Omega_T] = y_t y'_t \quad \text{and} \quad \mathbb{E}_{\theta(r)} [y_t f'_t | \Omega_T] = y_t \mathbb{E}_{\theta(r)} [f'_t | \Omega_T] \quad (10)$$

and $\mathbb{E}_{\theta(r)} [f_t | \Omega_T]$, $\mathbb{E}_{\theta(r)} [f_t f'_t | \Omega_T]$ and $\mathbb{E}_{\theta(r)} [f_t f'_{t-1} | \Omega_T]$ can be obtained through the Kalman smoother for the state space representation (4) (with the parameters $\theta(r)$). Shumway and Stoffer (1982) provide the estimates also for the incomplete data set however in the case when Λ is known. We provide the EM steps for the general case when Λ is unknown and y_t contains some missing values. In that case (10) no longer holds. The formulas (6) and (8) remain unaffected, however (7) and (9) need to be modified as follows:

$$\text{vec}(\Lambda(r+1)) = \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t f'_t | \Omega_T] \otimes W_t \right)^{-1} \text{vec} \left(\sum_{t=1}^T W_t y_t \mathbb{E}_{\theta(r)} [f'_t | \Omega_T] \right)$$

and

$$R(r+1) = \text{diag} \left(\frac{1}{T} \sum_{t=1}^T \left(W_t y_t y'_t W'_t - W_t y_t \mathbb{E}_{\theta(r)} [f'_t | \Omega_T] \Lambda(r+1)' W_t - W_t \Lambda(r+1) \mathbb{E}_{\theta(r)} [f_t | \Omega_T] y'_t W'_t \right. \right. \\ \left. \left. + W_t \Lambda(r+1) \mathbb{E}_{\theta(r)} [f_t f'_t | \Omega_T] \Lambda(r+1)' W_t + (I - W_t) R(r) (I - W_t) \right) \right),$$

where W_t is a diagonal matrix of size n with i^{th} diagonal element equal to 0 if $y_{i,t}$ is missing and equal to 1 otherwise. The derivations of the formulas are provided in the Appendix along with the results of simulations similar to those in Doz, Giannone, and Reichlin (2006) that verify the convergence of the estimates for different portions of missing data.

2.2 Forecasting

Given the estimates of the parameters $\hat{\theta} = (\hat{\Lambda}, \hat{A}_1, \hat{R}, \hat{Q})$ and the data set Ω_T the best linear predictions: $\mathbb{P}(y_{T+h}|\Omega_T)$, $h \geq 0$ can be obtained from the Kalman filter. In case some of the observations in y_t are missing, the corresponding rows in y_t and $\hat{\Lambda}$ (and the corresponding rows and columns in \hat{R}) are skipped when applying the Kalman filter (cf. Durbin and Koopman, 2001). The interpolates and backdata ($h < 0$) can be also simply obtained from the Kalman smoother.

To deal with mixed frequencies, low-frequency series are treated as high-frequency variables with missing data.

2.3 News and forecast revisions

In this section we show how the unexpected content, i.e. the *news*, in a data release is linked to the resulting forecast revision.

Let Ω_{v-1} and Ω_v be two consecutive vintages of data, consequently $\Omega_{v-1} \subset \Omega_v$.⁸ Let I_v denote the *news* in Ω_v with respect to Ω_{v-1} . For example, let us assume that the difference between Ω_{v-1} and Ω_v is the release of the European Commission (EC) Surveys for the period t_i . The *news* is $I_v = y_{t_i}^{EC} - \mathbb{P}(y_{t_i}^{EC}|\Omega_{v-1})$, where $y_{t_i}^{EC}$ is the vector of the new observations.

Assume that we are interested in how this *news* revises the GDP forecast for the period t_j . As $I_v \perp \Omega_{v-1}$ we can write

$$\mathbb{P}(\cdot|\Omega_v) = \mathbb{P}(\cdot|\Omega_{v-1}) + \mathbb{P}(\cdot|I_v)$$

or

$$\underbrace{\mathbb{P}(y_{t_j}^{GDP}|\Omega_v)}_{\text{new forecast}} = \underbrace{\mathbb{P}(y_{t_j}^{GDP}|\Omega_{v-1})}_{\text{old forecast}} + \underbrace{\mathbb{P}(y_{t_j}^{GDP}|I_v)}_{\text{news}}.$$

In other words, the updated forecast can be decomposed into the sum of the old forecast and of the contribution from the *news* in the latest release.

To compute the latter we use the fact that

$$\mathbb{P}(y_{t_j}^{GDP}|I_v) = \mathbb{E}(y_{t_j}^{GDP} I_v') \mathbb{E}(I_v I_v')^{-1} I_v.$$

Furthermore, given the model (4) we can write

$$\begin{aligned} y_{t_j}^{GDP} &= \lambda_{GDP} f_{t_j} + \xi_{t_j}^{GDP}, \\ I_v &= y_{t_j}^{EC} - y_{t_j}^{EC}|\Omega_{v-1} = \lambda_{EC} (f_{t_j} - f_{t_j}|\Omega_{v-1}) + \xi_{t_j}^{EC}, \end{aligned}$$

⁸In what follows, we do not take into account the revisions and changes in the parameter estimates. The influence of those factors needs to be analysed separately.

where λ_{GDP} and λ_{EC} are the rows of Λ corresponding to GDP and EC Surveys, respectively. It can be shown (see the Appendix) that:

$$\begin{aligned} E\left(y_{t_j}^{GDP} I_v'\right) &= \lambda_{GDP} E(f_{t_j} - f_{t_j|\Omega_{v-1}})(f_{t_i} - f_{t_i|\Omega_{v-1}})' \lambda_{EC}' \quad \text{and} \\ E(I_v I_v') &= \lambda_{EC} E(f_{t_i} - f_{t_i|\Omega_{v-1}})(f_{t_i} - f_{t_i|\Omega_{v-1}})' \lambda_{EC}' + R_{EC}, \end{aligned}$$

where R_{EC} is a diagonal matrix with elements of R corresponding to the EC Surveys. The expectations $E(f_{t_j} - f_{t_j|\Omega_{v-1}})(f_{t_i} - f_{t_i|\Omega_{v-1}})'$ and $E(f_{t_i} - f_{t_i|\Omega_{v-1}})(f_{t_i} - f_{t_i|\Omega_{v-1}})'$ can be obtained from the Kalman filter, see the Appendix.

Consequently, we can find a vector B such that the following holds:

$$\underbrace{y_{t_j|\Omega_v}^{GDP}}_{\text{new forecast}} = \underbrace{y_{t_j|\Omega_{v-1}}^{GDP}}_{\text{old forecast}} + B \underbrace{\left(y_{t_i}^{EC} - y_{t_i|\Omega_{v-1}}^{EC}\right)}_{\text{news}}. \quad (11)$$

This enables us trace the sources of forecast revisions.⁹ More precisely, in the case of a simultaneous release of several (groups of) variables it is possible to decompose the resulting forecast revision into contributions from the *news* in individual (groups of) series, see the illustration in Section 3.3.¹⁰ In addition, we can produce statements like e.g. “after the release of the EC Surveys, the forecast of GDP went up because the indicators turned out to be (on average) higher than expected”.¹¹

3 Forecasting the euro area GDP

In the empirical part of the paper we apply the methodology described above in the context of short-term forecasting and backcasting of the euro area GDP. In particular, we analyse the role of quarterly variables and the Purchasing Managers’ Surveys.

3.1 Data

The dataset contains in total 114 indicators and was downloaded in September 2007 directly after the release of industrial production series. The detailed description including list of the series, their availability and applied transformations is provided in the Appendix.

For the purpose of the empirical exercise the dataset is divided into three categories. The first group consists of 70 indicators similar to the one described in Bańbura and Rünstler (2007)

⁹Note, that the contribution from the *news* is equivalent to the change in the overall contribution of the series to the forecast (the measure proposed in Bańbura and Rünstler, 2007) when the correlations between the predictors are not exploited in the model. Otherwise, those measures are different. In particular, there can be a change in the overall contribution of a variable even if no new information on this variable was released.

¹⁰If the release concerns only one group or one series, the contribution of its *news* is simply equal to the change in the forecast.

¹¹This holds of course for the indicators with positive entries in B .

or used at the ECB for the short term forecast, cf. ECB (2008). It contains information on world prices, trade, industrial production, European Commission Surveys (EC Surveys), retail trade, labour, financial, US and interest rates. All the series in this group date back at least to January 1993, but most of them have a longer history. We refer to this group as the *Benchmark* dataset.

The second group contains 16 monthly “short-history” indicators, namely the Purchasing Managers’ Surveys (PM Surveys). They are available from August 1997 or July 1998. We refer to this group of data as *Short Monthly* dataset.

Finally the *Quarterly* dataset includes quarterly variables (with different time spans) on employment, unit labour cost, hourly labour cost, real GDP, real value added and EC Survey on capacity utilisation.

Following Giannone, Reichlin, and Small (2008) the monthly indicators are transformed so as to be in line with the nature of the quarterly variables (which are differences of 3 month totals).

3.2 Forecast evaluation

Given the trade-off between the timeliness and the information content, in the evaluation we want to account for the differences in the publication delays. In this, we aim at replicating as closely as possible the real-time forecasting exercise. As the real-time vintages are not available for all the variables of interest and whole evaluation period, we perform the so called pseudo-real time exercise. It means that we use the final figures as of mid-September 2007 however at each point of the evaluation sample we apply appropriate publication lags following the availability pattern for mid-September (this relies on the assumption that the release dates do not change much from month to month). For example, in mid-September the last available figure on the industrial production is for July. Consequently when we evaluate the model in e.g. April the data for the industrial production ”ends” in February. The same mechanism is applied to all the variables. The procedure for quarterly variables follows a similar logic modified to take into account the quarterly frequency of the releases.

To deal with mixed frequency data it is assumed that y_t are observed at a monthly frequency and the quarterly indicators are to be understood as monthly series with missing values in the first and second month of each quarter. Consequently, the GDP figure for a given quarter is assigned to the last month of this quarter.

We evaluate forecast accuracy for different forecast horizons: 0- (nowcast), 1- and 2-quarters ahead. In the case of nowcast we project the current quarter. This is relevant as the euro area GDP is released only around 6 weeks after the end of the respective quarter. Moreover, the available information changes depending on the respective month within a given quarter.

Therefore we evaluate separately forecasts made in the first, second or third month of a quarter. For example, 1-quarter-ahead forecast made in the second month means that we make forecast for the second quarter relying on the information available in February¹², for the third quarter using the information available in May, etc.

All the results in the paper correspond to a model with two factors ($r = 2$).¹³

The evaluation period is eight years, going from 1999Q3 to 2007Q2. For the measure of prediction accuracy we choose the mean squared forecast error (MSFE). All the tables present MSFEs relative to the a naïve constant growth model.¹⁴

The estimation sample starts in July 1993. We choose a recursive estimation which means that the sample length increases each time that more information becomes available.

We start by comparing the maximum likelihood method (ML) with the two-step approach (2S) proposed by Doz, Giannone, and Reichlin (2005) and used currently at the ECB for the short-term forecasting. One of the differences between the two methodologies applied on the benchmark dataset is in the way in which GDP is included in the analysis. In the 2S approach, the GDP is projected on the factors bridged to quarterly frequency and, given the parameters and the factor forecast produced by the factor model, the forecast of the GDP growth is computed. In the ML case the GDP is added to the benchmark dataset. This means that the GDP series will be the only one with missing values not only at the end of sample, but also in the first two months of each quarter.

Table 1: MSFEs for maximum likelihood and two-step approach

month	Nowcast		1-quarter-ahead		2-quarter-ahead	
	2S	ML	2S	ML	2S	ML
1	0.63	0.64	0.83	0.83	0.94	0.96
2	0.60	0.61	0.74	0.75	0.89	0.93
3	0.59	0.58	0.68	0.68	0.89	0.90

Notes: comparison between the MSFEs ratios to the constant growth rate model of the nowcast, 1-quarter-ahead and 2-quarter-ahead forecasts produced with a factor model estimated on the benchmark dataset with the maximum likelihood (ML) and the two-step (2S) approach.

Table 1 reports the corresponding MSFEs for the two methods. We can see that both produce comparable results. ML, however, has the advantage that it is easier to implement as it

¹²i.e. second month of first quarter;

¹³The results are qualitatively the same for smaller/larger number of factors.

¹⁴Let X_t and x_t denote the log of GDP and its quarterly growth rate, respectively (observed at quarterly frequency). The naïve constant growth model implies that X_t is a random walk with drift: $X_{t+1} = \delta + X_t + \epsilon_{t+1}$. Consequently $x_{t+1} = \delta + \epsilon_{t+1}$ and the optimal predictor is the average of the past growth rates and it does not depend on the forecast horizon h :

$$\hat{x}_{t+h|t} = \hat{\delta} = \frac{1}{t - T_0 + 1} \sum_{k=T_0}^t \hat{x}_k.$$

produces automatically the forecast for the GDP without the need for the second step.

In the following exercise we want to explore the effects of adding data with shorter history or quarterly frequency on the forecast accuracy. Table 2 presents the results for the nowcast and 1- and 2-quarter-ahead forecasts, respectively. In all the cases the parameters are estimated by maximum likelihood, what differs is the composition of the datasets. The first column reports the results obtained on the benchmark dataset (Bench); the second column reports the results obtained on the benchmark dataset augmented by the short monthly PM Surveys data (B+Monthly); the third column corresponds to the benchmark dataset plus all the quarterly variables (B+Quarterly); finally the fourth column reports the results obtained using all the data together (All). Figure 1 shows the true GDP growth and the nowcasts obtained with the four different datasets. The naïve “constant growth” projection is also plotted.

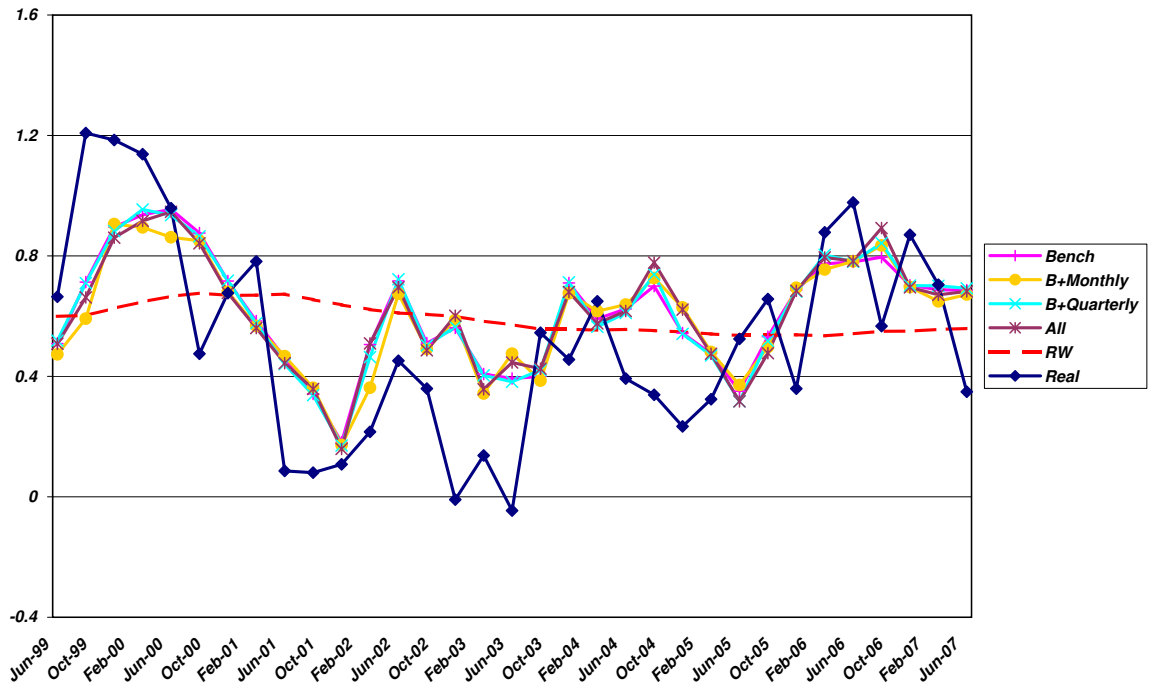
Table 2: MSFEs for maximum likelihood approach with different datasets

month	Bench	B+Monthly	B+Quarterly	All
Nowcast				
1	0.64	0.67	0.63	0.67
2	0.61	0.63	0.60	0.62
3	0.58	0.61	0.58	0.65
1-quarter-ahead forecast				
1	0.83	0.85	0.83	0.85
2	0.75	0.77	0.75	0.77
3	0.68	0.70	0.67	0.71
2-quarter-ahead forecast				
1	0.96	1.07	0.96	1.06
2	0.93	0.98	0.92	0.97
3	0.90	0.90	0.90	0.89

Notes: MSFEs ratios with respect to the constant growth rate model of the nowcast and forecasts produced with a factor model estimated on the benchmark dataset (Bench), on the benchmark dataset plus the short monthly indicators, on the benchmark dataset plus the quarterly indicators (B+Quarterly) and on the benchmark dataset plus the short monthly and the quarterly indicators (All).

We can see that the accuracy of the projections based on different datasets is comparable. Moreover, as Figure 1 shows for the nowcast, the differences between the point estimates are very small. In other words, there are no gains in forecast accuracy from augmenting the dataset by the additional indicators. On the other hand, the accuracy of forecasts based on the extended datasets does not deteriorate (perhaps with the exception of 2-quarter-ahead forecast). Consequently, we can include the additional series in the dataset and obtain their forecasts (or back estimates). This also allows us to analyse the *news* they provide, see the next section.

Figure 1: Nowcast of the GDP



3.3 News and forecast revisions, illustration

Let us recall from Section 2.3 that the *news* is understood as the unexpected part of a data release.

In the following exercise, on the example of the fourth quarter of 2001, we illustrate how the projection evolves with new data releases and what are the contributions of the *news* in different groups of series to the forecast revisions. We use the benchmark dataset augmented by the PM Surveys and we divide the indicators into 4 groups: industrial production indicators (IP), EC and PM Surveys and Other. We start to forecast in June 2001 (corresponding to the 2-quarter-ahead forecast in the third month) and we revise the projection each month as new data arrive. The last estimate is obtained in February 2002 and the actual GDP for the fourth quarter is released in March. Note that the 2 last projections are actually “backcasts” - they refer to the previous quarter. At each step we break down the forecast revision into the contributions of the *news* from different data groups using the formula (11).¹⁵

Figure 2 shows the evolution of the forecast as the new information arrives, the actual value of the GDP for the fourth quarter and the decomposition of the revisions.¹⁶ On average, the biggest contributions to the revisions come from the EC Surveys. In contrast, those from the industrial production have a sizeable effect on the projection only in the case of the backcast. This confirms the results of Bańbura and Rünstler (2007) on the important role of soft data for the GDP projections when the hard data for the relevant periods are not yet available.

Another observation is that the EC and PM Surveys do not always carry the same information. There are cases when their contributions are of the same sign and the cases when the opposite holds. In particular, in the case of nowcast the *news* in the EC Surveys “moves” the projection “towards” the actual value. This is not always the case for the PM Surveys.

3.4 Backdating the GDP

A useful feature of our framework is that the estimates of the missing observations in the panel can be obtained from the Kalman smoother in an automatic manner. This enables us to use it e.g. in order to backdate a short history series or to interpolate quarterly variables. In this section we evaluate the model at backdating the GDP. We include the GDP quarterly growth rate only as of June 2000 and discard all the previous observations. We run the model on the four datasets used in the forecast evaluation. Figure 3 plots the back estimates of the GDP from different datasets, the estimate obtained with a constant growth rate model and the actual quarterly growth rate of the GDP. Table 3 reports the mean squared errors of the

¹⁵The contribution of a group of series is the sum of the contributions of the series within this group.

¹⁶In fact “Other” contains also the effect of re-estimation.

Figure 2: Contribution of news to forecast revisions for 2001Q4

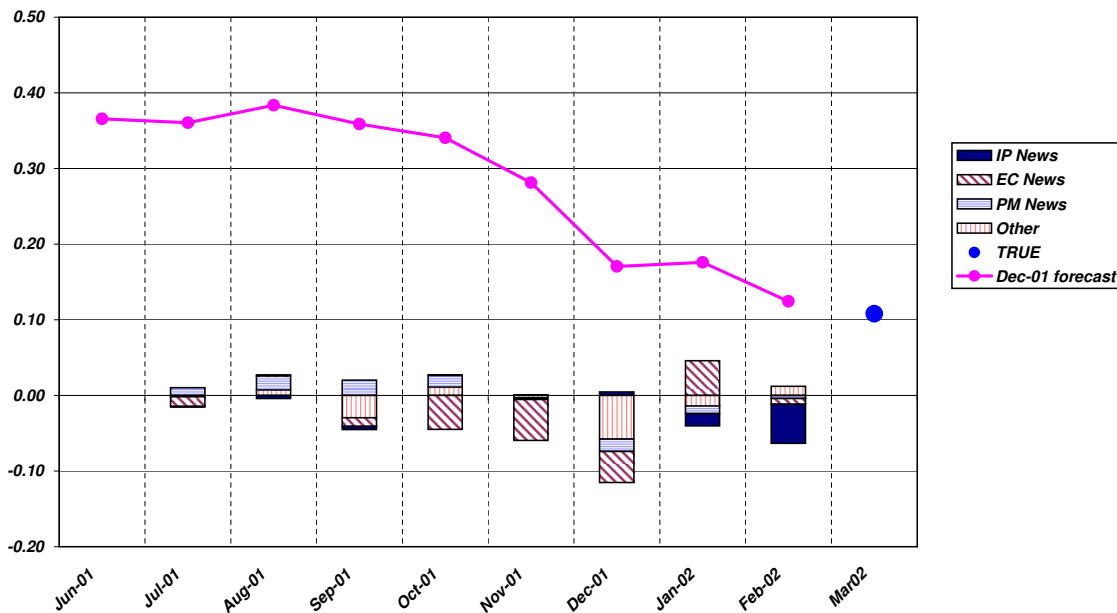
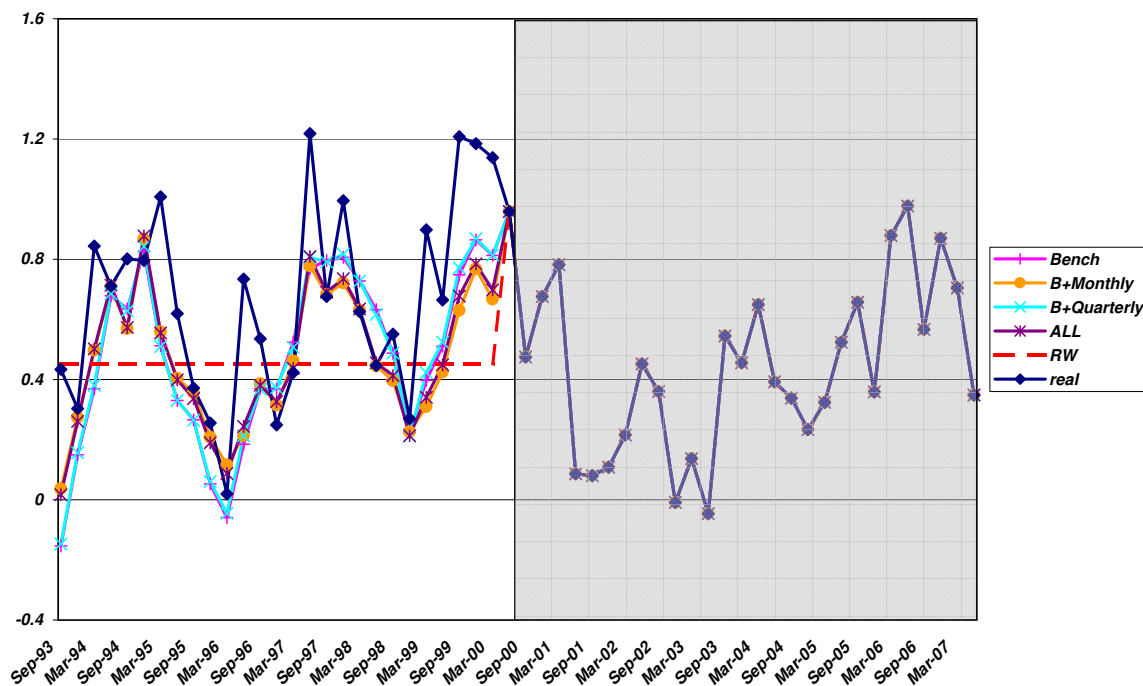


Figure 3: Back estimates of the GDP



backdated GDP obtained with the four different datasets relative to the naive constant growth rate estimate.

Table 3: MSEs of back estimates of the GDP

Bench	B+Monthly	B+Quarterly	ALL
0.59	0.58	0.55	0.53

Notes: MSEs relative to the constant growth rate model of the backcast produced with a factor model estimated on the benchmark dataset (Bench) including the “trimmed” GDP, on the benchmark dataset augmented by the short monthly indicators (B+Monthly), on the benchmark dataset plus the quarterly indicators (B+Quarterly) and on the entire dataset (All).

As we can see from Table 3 and Figure 3, independently of the dataset used, the backdating seems to capture the movements of the GDP, giving reasonable estimates of the past values of the series. Again, different datasets yield comparable results.

4 Summary

This paper proposes a methodology for the estimation of a large dynamic factor model for datasets with arbitrary pattern of missing values. For that we adopt the maximum likelihood approach, which is shown by Doz, Giannone, and Reichlin (2005) to be consistent as well as computationally feasible for large cross-sections once the Expectation-Maximisation (EM) algorithm is applied. We show how the steps of the EM algorithm should be modified in the case of missing data. In addition, we derive the link between the unexpected component of data releases and the resulting forecast revision.

We apply this methodology for the short-term forecasting of the euro area GDP on the basis of large monthly dataset. Thanks to the flexibility of the framework in dealing with missing data, short history and quarterly variables can be also considered (e.g. Purchasing Managers' Surveys, GDP components or labour statistics). The effect of including these indicators in the large monthly dataset similar to the one used in e.g. Bańbura and Rünstler (2007) or ECB (2008) is evaluated in an out-of-sample forecast exercise. The results indicate that the additional indicators do not improve the precision of the projections. On the other hand, they can be analysed and forecast in a single model. Finally, we show that the framework can be easily used for back estimation. In particular, the back estimates of the GDP are fairly close to the true values.

References

- AASTVEIT, K., AND T. TROVIK (2007): "Nowcasting Norwegian GDP: The role of asset prices in a small open economy," Norges Bank Working Paper 2007/09.
- ANGELINI, E., J. HENRY, AND M. MARCELLINO (2006): "Interpolation and backdating with a large information set," *Journal of Economic Dynamics and Control*, 30(12), 2693–2724.
- BAÑBURA, M., AND G. RÜNSTLER (2007): "A look into the factor model black box. Publication lags and the role of hard and soft data in forecasting GDP.," Working Paper Series 751, European Central Bank.
- BERNANKE, B. S., AND J. BOIVIN (2003): "Monetary policy in a data-rich environment," *Journal of Monetary Economics*, 50(3), 525–546.
- BOIVIN, J., AND S. NG (2005): "Understanding and Comparing Factor-Based Forecasts," *International Journal of Central Banking*, 3, 117–151.
- BREITUNG, J., AND C. SCHUMACHER (2008): "Real-time forecasting of German GDP based on a large factor model with monthly and quarterly Data," *International Journal of Forecasting*, 24, 386–398.

- CAMACHO, M., AND G. PEREZ-QUIROS (2008): “Introducing the EURO-STING: Short Term Indicator of Euro Area Growth,” Banco de España Working Papers 0807, Banco de España.
- D’AGOSTINO, A., AND D. GIANNONE (2006): “Comparing alternative predictors based on large-panel factor models,” Working Paper Series 680, European Central Bank.
- DEMPSTER, A., N. LAIRD, AND D. RUBIN (1977): “Maximum Likelihood Estimation From Incomplete Data,” *Journal of the Royal Statistical Society*, 14, 1–38.
- DOZ, C., D. GIANNONE, AND L. REICHLIN (2005): “A two-step estimator for large approximate dynamic factor models based on Kalman filtering,” Manuscript, Université Libre de Bruxelles.
- (2006): “A Quasi Maximum Likelihood Approach for Large Approximate Dynamic Factor Models,” Working Paper Series 674, European Central Bank.
- DURBIN, J., AND S. J. KOOPMAN (2001): *Time Series Analysis by State Space Methods*. Oxford University Press.
- ECB (2008): “Short-term forecasts of economic activity in the euro area,” in *Monthly Bulletin*, April, pp. 69–74. European Central Bank.
- FORNI, M., M. HALLIN, M. LIPPI, AND L. REICHLIN (2000): “The Generalized Dynamic Factor Model: identification and estimation,” *Review of Economics and Statistics*, 82, 540–554.
- (2003): “Do Financial Variables Help Forecasting Inflation and Real Activity in the Euro Area?,” *Journal of Monetary Economics*, 50, 1243–55.
- (2005): “The Generalized Dynamic Factor Model: one-sided estimation and forecasting,” *Journal of the American Statistical Association*, 100, 830–840.
- GEWEKE, J. F. (1977): “The Dynamic Factor Analysis of Economic Time Series Models,” in *Latent Variables in Socioeconomic Models*, ed. by D. Aigner, and A. Goldberger, pp. 365–383. North-Holland.
- GIANNONE, D., L. REICHLIN, AND L. SALA (2004): “Monetary Policy in Real Time,” in *NBER Macroeconomics Annual*, ed. by M. Gertler, and K. Rogoff, pp. 161–200. MIT Press.
- GIANNONE, D., L. REICHLIN, AND D. SMALL (2008): “Nowcasting: The real-time informational content of macroeconomic data,” *Journal of Monetary Economics*, 55(4), 665–676.
- JUNGBACKER, B., AND S. J. KOOPMAN (2008): “Likelihood-based Analysis for Dynamic Factor Models,” Tinbergen Institute Discussion Papers 08-007/4, Tinbergen Institute.

- KOOPMAN, S. J., AND J. DURBIN (2000): “Fast filtering and smoothing for multivariate state space models,” *Journal of Time Series Analysis*, 21(3), 281–296.
- MARCELLINO, M., J. H. STOCK, AND M. W. WATSON (2003): “Macroeconomic forecasting in the Euro area: Country specific versus area-wide information,” *European Economic Review*, 47(1), 1–18.
- MATHESON, T. (2007): “An analysis of the informational content of New Zealand data releases: the importance of business opinion surveys,” Reserve Bank of New Zealand Discussion Paper 2007/13.
- MCLACHLAN, G. J., AND T. KRISHNAN (1996): *The EM Algorithm and Extensions*. John Wiley and Sons.
- PROIETTI, T. (2008): “Estimation of Common Factors under Cross-Sectional and Temporal Aggregation Constraints: Nowcasting Monthly GDP and its Main Components,” Manuscript.
- REIS, R., AND M. W. WATSON (2007): “Relative Goods’ Prices and Pure Inflation,” NBER Working Paper 13615.
- SARGENT, T. J., AND C. SIMS (1977): “Business Cycle Modelling without Pretending to have too much a-priori Economic Theory,” in *New Methods in Business Cycle Research*, ed. by C. Sims. Federal Reserve Bank of Minneapolis.
- SHUMWAY, R., AND D. STOFFER (1982): “An approach to time series smoothing and forecasting using the EM algorithm,” *Journal of Time Series Analysis*, 3, 253–264.
- STOCK, J. H., AND M. W. WATSON (2002a): “Forecasting Using Principal Components from a Large Number of Predictors,” *Journal of the American Statistical Association*, 97, 1167–1179.
- (2002b): “Macroeconomic Forecasting Using Diffusion Indexes.,” *Journal of Business and Economics Statistics*, 20, 147–162.
- WATSON, M. W., AND R. F. ENGLE (1983): “Alternative algorithms for the estimation of dynamic factor, mimic and varying coefficient regression models,” *Journal of Econometrics*, 23, 385–400.

A Data Description

Number	Block	Name	Transf.	Availab.
<i>Benchmark dataset</i>				
1	World Prices	World market prices of raw materials in Euro. Index Total.	3	Jan-80
2		World market prices of raw materials in Euro. Index Total, excluding Energy	3	Jan-80
3		WORLD-MKT PRICES, ENERGY RAW MAT., CRUDE OIL	3	Jan-80
4		GOLD PRICE, US DOLLARS/FINE OUNCE,LONDON FIXING	3	Jan-80
5		Brent Crude-1 Month Fwd,fob US\$/BBL converted in euro	3	May-85
6	Trade	trade with World (all entities), Export Value SA, not w.d. adj.	2	Jan-80
7		trade with World (all entities), Export Value, SA, not w.d. adj.	2	Jan-80
8		trade with World (all entities), Export Value, SA, not w.d. adj.	2	Jan-80
9		trade with World (all entities), Export Value, SA, not w.d. adj.	2	Jan-80
10	Industrial Production	Retail trade, except of motor vehicles and motorcycles - W.d. and SA	2	Jan-80
11		Total industry - w.d. and SA	2	Jan-88
12		Total Industry (excluding construction) - W.d. and SA	2	Jan-85
13		Manufacturing - W.d. and SA	2	Jan-85
14		Construction - W.d. and SA	2	Jan-88
15		Total Industry excluding Construction and MIG Energy - W.d. and SA	2	Jan-85
16		Energy excluding NACE Rev.1 Section E - W.d. and SA	2	Jan-90
17		MIG Capital Goods Industry - W.d. and SA	2	Jan-80
18		MIG Durable Consumer Goods Industry - W.d. and SA	2	Jan-90
19		MIG Energy - Working day and SA	2	Jan-80
20		MIG Intermediate Goods Industry - Working day and SA	2	Jan-80
21		MIG Non-durable Consumer Goods Industry - W.d. and SA	2	Jan-85
22		Manufacture of basic metals - W.d. and SA	2	Jan-85
23		Manufacture of chemicals and chemical products - W.d. and SA	2	Jan-85
24		Manufacture of electrical machinery and apparatus n.e.c. - W.d. and SA	2	Jan-80
25		Manufacture of machinery and equipment n.e.c. - W.d. and SA	2	Jan-85
26		Manufacture of pulp, paper and paper products - W.d. and SA	2	Jan-85
27		Manufacture of rubber and plastic products - W.d. and SA	2	Jan-85
28	European Commission Surveys	Industry Survey: Industrial Confidence Indicator - SA	1	Jan-85
29		Industry Survey: Production trend observed in recent months - SA	1	Jan-85
30		Industry Survey: Assessment of order-book levels - SA	1	Jan-85
31		Industry Survey: Assessment of export order-book levels - SA	1	Jan-85
32		Industry Survey: Assessment of stocks of finished products - SA	1	Jan-85
33		Industry Survey: Production expectations for the months ahead - SA	1	Jan-85
34		Industry Survey: Employment expectations for the months ahead - SA	-	Jan-85
35		Consumer Survey: Consumer Confidence Indicator - SA	1	Jan-85
36		Consumer Survey: General economic situation over last 12 months - SA	1	Jan-85
37		Consumer Survey: General economic situation over next 12 months - SA	1	Jan-85
38		Consumer Survey: Unemployment expectations over next 12 months - SA	1	Jan-85
39		Construction Survey: Construction Confidence Indicator - SA	1	Jan-85
40		Construction Survey: Trend of activity compared with preceding months 1 SA	-	Jan-85
41		Construction Survey: Assessment of order books - SA	1	Jan-85
42		Construction Survey: Employment expectations for the months ahead - SA	1	Jan-85
43	Retail Trade	Retail Trade Survey: Retail Confidence Indicator - SA	1	Jan-85
44		Retail Trade Survey: Present business situation - SA	1	Jan-85
45		Retail Trade Survey: Assessment of stocks - SA	1	Jan-85
46		Retail Trade Survey: Expected business situation - SA	1	Jan-85
47		Retail Trade Survey: Employment expectations - SA	1	Apr-85
48	New passenger car - W.d. and SA	2	Jan-90	
49	Labour	Unemployment rate, Total (all ages), Total (male & female), SA, not w.d. adj	1	Jan-93
50		Index of Employment, Construction; SA, not w.d. adj	2	Jan-93
51		Index of Employment, Manufacturing; SA, not w.d. adj	2	Oct-89
52		Index of Employment, Total Industry; SA, not w.d. adj	2	Jan-93
53		Index of Employment, Total Industry (excluding construction);	2	Oct-89
54	Financial	EUROSTOXX 50 (RHS)	2	Dec-86
55		EUROSTOXX 325 (LHS)	2	Dec-86
56		S&P 500 COMPOSITE - PRICE INDEX	2	Jan-80
57	US	US, STOCK-EXCH. PRICES, DOW JONES, INDUSTRIAL AVERAGE, NSA	2	Jan-80
58		US, INT.RATE, MONEY-MKT, TREAS.BILLS, 3-MONTH, MKT YIELD	1	Jan-80
59		US, YIELD, SECOND.MKT, US TREASURY NOTES & BONDS, 10 YEARS	1	Jan-80
60		US, UNEMPLOYMENT RATE, SA	1	Jan-80
61		US, INDUST.PROD., TOTAL EXCL. CONSTRUCTION, SA	2	Jan-80
62		US, EMPLOYMENT, CIVILIAN, SA	2	Jan-80
63		US, RETAIL TRADE, VALUE (NAICS DEF.), SA	2	Jan-92
64		US, PRODUCTION EXPECTATIONS IN MANUF., NSA	1	Jan-80
65		US, CONSUMER EXPECTATIONS INDEX, NSA	1	Jan-80
66	Interest Rates	10-year govt. bonds	1	Jan-80
67		Reuters.Money market.Euro.Euribor	1	Jan-80
68		1-year EURIBOR RATE	1	Jan-80
69		YIELD, SECOND MKT, GOVT BONDS, 2 YEARS	1	Jan-80
70		YIELD, SECOND MKT, GOVT BONDS, 5 YEARS	1	Jan-80
<i>Short Monthly dataset</i>				
71	Purchasing Managers' Surveys	Manufacturing - employment	1	Aug-97
72		Manufacturing - new orders	1	Aug-97
73		Manufacturing - new export orders	1	Aug-97
74		Manufacturing - output	1	Aug-97
75		Manufacturing - purchasing manager index	1	Aug-97
76		Manufacturing - productivity	1	Jan-98
77		Manufacturing - quantity of purchases	1	Aug-97
78		Manufacturing - supplier delivery times	1	Aug-97
79		Manufacturing - stocks of finished goods	1	Aug-97
80		Manufacturing - stocks of purchases	1	Aug-97
81		Services - employment	1	Jul-98
82		Services - future business activity expectations	1	Jul-98
83		Services - new business	1	Jul-98
84		Services - outstanding business	1	Jul-98
85		Services - business activity	1	Jul-98
86		Services - productivity	1	Jan-98
<i>Quarterly dataset</i>				

87		Total domestic- Level - Thousands of persons - SA	4	Q1-1980
88		Employees: total domestic- Level - Thousands of persons - NSA	4	Q1-1980
89	Employment:	Self-employed: total- Level - Thousands of persons - NSA	4	Q1-1980
90	by employment status,	Total industry- Level - Thousands of persons - NSA	4	Q1-1980
91	by economic activity	Construction- Level - Thousands of persons - NSA	4	Q1-1980
92		Trade and other - Level - Thousands of persons - NSA	4	Q1-1990
93		Total - index 2000 = 100 - SA	4	Q1-1995
94	Unit labour cost:	Industry, including energy - Index 2000 = 100 - SA	4	Q1-1995
95	by economic activity	Construction - Index 2000 = 100 - SA	4	Q1-1995
96		Trade and other - Index 2000 = 100 - SA	4	Q1-1995
97	Hourly labour cost	Whole economy excluding agriculture, fishing and government, SA	4	Q1-1996
98	US	UNIT LABOUR COSTS IN MANUFACTURING, SA	4	Q1-1959
99		GDP, AT MARKET PRICES - CHAINED 2000 USD SAAR	4	Q1-1959
100		Gross domestic product at market price - Chain linked - SA	4	Q1-1991
101		Final consumption of households and NPISH's - Chain linked - SA	4	Q1-1991
102	GDP: by expenditure	Final consumption of general government - Chain linked - SA	4	Q1-1991
103	components at	Gross fixed capital formation - Chain linked - SA	4	Q1-1991
104	constant prices	Exports of goods and services - Chain linked - SA	4	Q1-1991
105		Imports of goods and services - Chain linked - SA	4	Q1-1991
106		Gross value added at basic prices - SA	4	Q1-1995
107		Agricultural, hunting, forestry and fishing products - SA	4	Q1-1995
108	Value added:	Total industry - SA	4	Q1-1995
109	by economic activity	Construction - SA	4	Q1-1995
110	at constant prices	Trade and other - SA	4	Q1-1995
111		Financial intermediation, real estate - SA	4	Q1-1995
112		Other services - SA	4	Q1-1995
113		Taxes less subsidies on products - SA	4	Q1-1995
114	EC Survey	Industry Survey: Current level of capacity utilization	5	Q1-1980

Transformation code:

1: $y_{it} = (1 + L + L^2)(1 - L^3)Y_{it}$; 2: $y_{it} = (1 + L + L^2)(1 - L^3) \log(Y_{it})$;

3: $y_{it} = (1 + L + L^2)(1 - L^3)(1 - L^{12}) \log(Y_{it})$; 4: $y_{it} = (1 - L^3) \log(Y_{it})$; 5: $y_{it} = (1 - L^3)(Y_{it})$

B Derivation of the EM iterations

Let us first derive the estimates for the case of $p = 1$. Let us recall that $\theta = (\Lambda, A_1, R, Q)$. For the model given by (4) the joint log-likelihood (for the observations and the latent factors) is given by:

$$\begin{aligned}
l(Y, F; \theta) &= -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (f_0 - \mu)' \Sigma^{-1} (f_0 - \mu) \\
&- \frac{T}{2} \log |Q| - \frac{1}{2} \sum_{t=1}^T (f_t - A_1 f_{t-1})' Q^{-1} (f_t - A_1 f_{t-1}) \\
&- \frac{T}{2} \log |R| - \frac{1}{2} \sum_{t=1}^T (y_t - \Lambda f_t)' R^{-1} (y_t - \Lambda f_t) \\
&= -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (f_0 - \mu)' \Sigma^{-1} (f_0 - \mu) \\
&- \frac{T}{2} \log |Q| - \frac{1}{2} \text{tr} \left[Q^{-1} \sum_{t=1}^T (f_t - A_1 f_{t-1})(f_t - A_1 f_{t-1})' \right] \\
&- \frac{T}{2} \log |R| - \frac{1}{2} \text{tr} \left[R^{-1} \sum_{t=1}^T (y_t - \Lambda f_t)(y_t - \Lambda f_t)' \right]
\end{aligned}$$

By differentiating $\mathbb{E}_{\theta(r)} [l(Y, F; \theta) | \Omega_T]$ with respect to A_1 we get

$$\begin{aligned}
\frac{\partial \mathbb{E}_{\theta(r)} [l(Y, F; \theta) | \Omega_T]}{\partial A_1} &= -\frac{1}{2} \frac{\partial \text{tr} \left\{ Q^{-1} \sum_{t=1}^T \mathbb{E}_{\theta(r)} [(f_t - A_1 f_{t-1})(f_t - A_1 f_{t-1})' | \Omega_T] \right\}}{\partial A_1} \\
&= -Q^{-1} \sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t f_{t-1}' | \Omega_T] + Q^{-1} A_1 \sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_{t-1} f_{t-1}' | \Omega_T],
\end{aligned}$$

and consequently

$$A_1(r+1) = \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_t f_{t-1}' | \Omega_T] \right) \left(\sum_{t=1}^T \mathbb{E}_{\theta(r)} [f_{t-1} f_{t-1}' | \Omega_T] \right)^{-1},$$

as provided in the main text. Similarly we can obtain $\Lambda(r+1)$, $Q(r+1)$, and $R(r+1)$, given by the formulas (7)-(9)

In case y_t contains missing values, (10) no longer holds and the formulas for $\Lambda(r+1)$ and $R(r+1)$ need to be modified. Let us differentiate $\mathbb{E}_{\theta(r)} [l(Y, F; \theta) | \Omega_T]$ with respect to Λ :

$$\frac{\partial \mathbb{E}_{\theta(r)} [l(Y, F; \theta) | \Omega_T]}{\partial \Lambda} = -\frac{1}{2} \frac{\partial \text{tr} \left\{ R^{-1} \sum_{t=1}^T \mathbb{E}_{\theta(r)} [(y_t - \Lambda f_t)(y_t - \Lambda f_t)' | \Omega_T] \right\}}{\partial \Lambda},$$

and let us have a closer look at $\mathbb{E}[(y_t - \Lambda f_t)(y_t - \Lambda f_t)' | \Omega_T]$ (to simplify the notation we skip the subscript $\theta(r)$).

Let

$$y_t = W_t y_t + (I - W_t) y_t = y_t^{(1)} + y_t^{(2)},$$

where W_t is a diagonal matrix with ones corresponding to the non-missing entries in y_t and 0 otherwise. ($y_t^{(1)}$ contains the non-missing observations at time t with 0 in place of the missing ones.)

We have:

$$\begin{aligned}
(y_t - \Lambda f_t)(y_t - \Lambda f_t)' &= \left(W_t(y_t - \Lambda f_t) + (I - W_t)(y_t - \Lambda f_t) \right) \left(W_t(y_t - \Lambda f_t) + (I - W_t)(y_t - \Lambda f_t) \right)' \\
&= W_t(y_t - \Lambda f_t)(y_t - \Lambda f_t)'W_t + (I - W_t)(y_t - \Lambda f_t)(y_t - \Lambda f_t)'(I - W_t) \\
&\quad + W_t(y_t - \Lambda f_t)(y_t - \Lambda f_t)'(I - W_t) + (I - W_t)(y_t - \Lambda f_t)(y_t - \Lambda f_t)'W_t.
\end{aligned}$$

By the law of iterated expectations:

$$E[(y_t - \Lambda f_t)(y_t - \Lambda f_t)'|\Omega_T] = E\left[E[(y_t - \Lambda f_t)(y_t - \Lambda f_t)'|F, \Omega_T]|\Omega_T\right].$$

As

$$\begin{aligned}
E[W_t(y_t - \Lambda f_t)(y_t - \Lambda f_t)'(I - W_t)|F, \Omega_T] &= 0, \\
E[(I - W_t)(y_t - \Lambda f_t)(y_t - \Lambda f_t)'(I - W_t)|F, \Omega_T] &= (I - W_t)R(r)(I - W_t)
\end{aligned}$$

and

$$\begin{aligned}
&E[W_t(y_t - \Lambda f_t)(y_t - \Lambda f_t)'W_t|\Omega_T] \\
&= W_t y_t y_t' W_t - W_t y_t E[f_t'|\Omega_T] \Lambda' W_t - W_t \Lambda E[f_t|\Omega_T] y_t' W_t + W_t \Lambda E[f_t f_t'|\Omega_T] \Lambda' W_t \\
&= y_t^{(1)} y_t^{(1)'} - y_t^{(1)} E[f_t'|\Omega_T] \Lambda' W_t - W_t \Lambda E[f_t|\Omega_T] y_t^{(1)'} + W_t \Lambda E[f_t f_t'|\Omega_T] \Lambda' W_t,
\end{aligned}$$

we get:

$$\begin{aligned}
&E[(y_t - \Lambda f_t)(y_t - \Lambda f_t)'|\Omega_T] = \\
&y_t^{(1)} y_t^{(1)'} - y_t^{(1)} E[f_t'|\Omega_T] \Lambda' W_t - W_t \Lambda E[f_t|\Omega_T] y_t^{(1)'} + W_t \Lambda E[f_t f_t'|\Omega_T] \Lambda' W_t + (I - W_t)R(r)(I - W_t).
\end{aligned}$$

Consequently:

$$\begin{aligned}
\frac{\partial \text{tr}\left\{R^{-1}E[(y_t - \Lambda f_t)(y_t - \Lambda f_t)'|\Omega_T]\right\}}{\partial \Lambda} &= -2W_t R^{-1} y_t^{(1)} E[f_t'|\Omega_T] + 2W_t R^{-1} W_t \Lambda E[f_t f_t'|\Omega_T] \\
&= -2R^{-1} y_t^{(1)} E[f_t'|\Omega_T] + 2R^{-1} W_t \Lambda E[f_t f_t'|\Omega_T].
\end{aligned}$$

From

$$\sum_{t=1}^T \frac{\partial \text{tr}\left\{R^{-1}E_{\theta(r)}[(y_t - \Lambda f_t)(y_t - \Lambda f_t)'|\Omega_T]\right\}}{\partial \Lambda} \Bigg|_{\Lambda=\Lambda(r+1)} = 0$$

follows

$$\sum_{t=1}^T y_t^{(1)} E_{\theta(r)}[f_t'|\Omega_T] = \sum_{t=1}^T W_t \Lambda(r+1) E_{\theta(r)}[f_t f_t'|\Omega_T].$$

Equivalently (as $\text{vec}(ABC) = (C' \otimes A)\text{vec}(B)$) we have

$$\text{vec}\left(\sum_{t=1}^T y_t^{(1)} E_{\theta(r)}[f_t'|\Omega_T]\right) = \left(\sum_{t=1}^T E_{\theta(r)}[f_t f_t'|\Omega_T] \otimes W_t\right) \text{vec}(\Lambda(r+1))$$

hence

$$\text{vec}(\Lambda(r+1)) = \left(\sum_{t=1}^T \mathbf{E}_{\theta(r)} [f_t f_t' | \Omega_T] \otimes W_t \right)^{-1} \text{vec} \left(\sum_{t=1}^T y_t^{(1)} \mathbf{E}_{\theta(r)} [f_t' | \Omega_T] \right).$$

In the similar fashion we obtain

$$\begin{aligned} R(r+1) &= \text{diag} \left(\frac{1}{T} \sum_{t=1}^T \left(y_t^{(1)} y_t^{(1)'} - y_t^{(1)} \mathbf{E}_{\theta(r)} [f_t' | \Omega_T] \Lambda(r+1)' W_t - W_t \Lambda(r+1) \mathbf{E}_{\theta(r)} [f_t | \Omega_T] y_t^{(1)'} \right. \right. \\ &\quad \left. \left. + W_t \Lambda(r+1) \mathbf{E}_{\theta(r)} [f_t f_t' | \Omega_T] \Lambda(r+1)' W_t + (I - W_t) R(r) (I - W_t) \right) \right). \end{aligned}$$

Let us now consider the case of $p > 1$. We can write the log-likelihood:

$$\begin{aligned} l(Y, F; \theta) &= -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\tilde{f}_0 - \mu)' \Sigma^{-1} (\tilde{f}_0 - \mu) \\ &\quad - \frac{r}{2} \log |Q| - \frac{1}{2} \text{tr} \left[Q^{-1} \sum_{t=1}^T (f_t - A \tilde{f}_{t-1}) (f_t - A \tilde{f}_{t-1})' \right] \\ &\quad - \frac{n}{2} \log |R| - \frac{1}{2} \text{tr} \left[R^{-1} \sum_{t=1}^T (y_t - \Lambda f_t) (y_t - \Lambda f_t)' \right], \end{aligned}$$

where $\tilde{f}_{t-1} = [f'_{t-1}, \dots, f'_{t-p}]'$ and $A = [A_1, \dots, A_p]$.

Consequently (6) and (8) should be modified as:

$$A(r+1) = \left(\sum_{t=1}^T \mathbf{E}_{\theta(r)} [f_t \tilde{f}'_{t-1} | \Omega_T] \right) \left(\sum_{t=1}^T \mathbf{E}_{\theta(r)} [\tilde{f}_{t-1} \tilde{f}'_{t-1} | \Omega_T] \right)^{-1}$$

and

$$Q(r+1) = \frac{1}{T} \left(\sum_{t=1}^T \mathbf{E}_{\theta(r)} [f_t f_t' | \Omega_T] - A(r+1) \sum_{t=1}^T \mathbf{E}_{\theta(r)} [\tilde{f}_{t-1} \tilde{f}'_{t-1} | \Omega_T] \right).$$

The conditional moments of the factors $\mathbf{E}_{\theta(r)} [f_t \tilde{f}'_{t-1} | \Omega_T]$, $\mathbf{E}_{\theta(r)} [\tilde{f}_{t-1} \tilde{f}'_{t-1} | \Omega_T]$, $\mathbf{E}_{\theta(r)} [f_t f_t' | \Omega_T]$ can be obtained by running the Kalman filter on the following state space form:

$$\begin{aligned} Y_t &= \begin{bmatrix} \Lambda & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} f_t \\ f_{t-1} \\ \vdots \\ f_{t-p+1} \end{bmatrix} + \epsilon_t \quad \epsilon_t \sim N(0, R), \\ \begin{bmatrix} f_t \\ f_{t-1} \\ \vdots \\ f_{t-p+1} \end{bmatrix} &= \begin{bmatrix} A_1 & A_2 & \dots & A_p \\ I & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & I & 0 \end{bmatrix} \begin{bmatrix} f_{t-1} \\ f_{t-2} \\ \vdots \\ f_{t-p} \end{bmatrix} + u_t \quad u_t \sim N \left(0, \begin{bmatrix} Q & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \right). \end{aligned}$$

C Convergence checks

We perform simulations similar to those in Doz, Giannone, and Reichlin (2006). More precisely we simulate the data from the following factor model:

$$\begin{aligned} y_t &= \Lambda f_t + \xi_t, \\ f_t &= A f_{t-1} + u_t \quad u_t \sim N(0, I_r), \\ \xi_t &= D \xi_{t-1} + v_t \quad v_t \sim N(0, \Phi), \end{aligned}$$

where

$$\begin{aligned} \Lambda_{ij} &\text{ i.i.d. } N(0, 1), \\ A_{ij} &= \begin{cases} \rho, & i = j \\ 0, & i \neq j \end{cases}, \quad D_{ij} = \begin{cases} d, & i = j \\ 0, & i \neq j \end{cases}, \\ \Phi &= \tau^{|i-j|} (1 - d^2) \sqrt{\alpha_i \alpha_j}, \quad \alpha_i = \frac{\beta_i}{1 - \beta_i} \frac{1}{1 - \rho^2} \sum_{j=1}^r \Lambda_{ij}^2, \quad \beta_i \text{ i.i.d. } U([u, 1 - u]), \end{aligned}$$

see Stock and Watson (2002a) or Doz, Giannone, and Reichlin (2006) for the discussion. This model allows for serial and (some degree of) cross-correlation in the idiosyncratic components.

The data are generated with $r = 3$, $\rho = 0.9$, $d = 0.5$, $u = 0.1$ and $\tau = 0.5$. Subsequently, we set $k\%$ of the data as missing and we estimate the parameters and the factors. We use the same stopping criterion as Doz, Giannone, and Reichlin (2006).

To assess the precision of the estimates of the factors we follow Stock and Watson (2002a) and Doz, Giannone, and Reichlin (2006) and use the trace R^2 of the regression of the estimated factors on the true ones:

$$\frac{\text{Trace}(F' \hat{F} (\hat{F}' \hat{F})^{-1} \hat{F}' F)}{\text{Trace}(F' F)}.$$

This measure is smaller than 1 and tends to 1 with the increasing canonical correlation between the estimated and the true factors.

Table 5 presents the trace statistics for different cross-section size n , different series length T and different missing-data ratio k .

We can see that in each case the factors converge to the true ones with increasing T and n albeit for larger share of missing data the convergence is slightly slower.

In the case with missing data the convergence of the EM algorithm is slower as well. To speed up the computations, the solutions proposed in Koopman and Durbin (2000) and Jungbacker and Koopman (2008) could be considered.¹⁷

D Computation of the news

As in the Section 2.3 let Ω_{v-1} and Ω_v be two consecutive vintages of data and let I_v be the *news* content of Ω_v orthogonal to Ω_{v-1} . We will derive the formula for the projection $\mathbb{P}(y_{j,t_j} | I_v)$ for a general form of I_v .

¹⁷In the case of the former the diagonal form of the covariance matrix of the idiosyncratic component can be exploited, in the case of the latter, the fact that the size of observation vector is large relative to the size of the state vector.

Table 5: Trace statistics for the factor estimates

T/n	10	25	50	100
No missing data				
50	0.54	0.64	0.69	0.72
100	0.61	0.76	0.81	0.83
10% missing data				
50	0.53	0.63	0.69	0.72
100	0.61	0.76	0.81	0.83
20% missing data				
50	0.53	0.63	0.68	0.71
100	0.60	0.75	0.80	0.83
40% missing data				
50	0.50	0.60	0.67	0.70
100	0.59	0.73	0.79	0.82

Notes: Table reports trace statistics for the factor estimates for different ratios of missing data for data simulated from a factor model. T and n refer to the sample and cross-section size respectively.

Without a loss of generality, let us assume that Ω_v contains new releases for the first K variables for the reference periods t_1, \dots, t_K , respectively. Then we have

$$I_v = \begin{bmatrix} y_{1,t_1} - \mathbb{P}(y_{1,t_1} | \Omega_{v-1}) \\ y_{2,t_2} - \mathbb{P}(y_{2,t_2} | \Omega_{v-1}) \\ \vdots \\ y_{K,t_K} - \mathbb{P}(y_{K,t_K} | \Omega_{v-1}) \end{bmatrix}.$$

The projection is given by

$$\mathbb{P}(y_{j,t_j} | I_v) = \mathbb{E}(y_{j,t_j} I_v') \mathbb{E}(I_v I_v')^{-1} I_v, \quad (12)$$

where

$$\mathbb{E}(y_{j,t_j} I_v') = \begin{bmatrix} \mathbb{E}(y_{j,t_j} (y_{1,t_1} - \mathbb{P}(y_{1,t_1} | \Omega_{v-1}))) \\ \mathbb{E}(y_{j,t_j} (y_{2,t_2} - \mathbb{P}(y_{2,t_2} | \Omega_{v-1}))) \\ \vdots \\ \mathbb{E}(y_{j,t_j} (y_{K,t_K} - \mathbb{P}(y_{K,t_K} | \Omega_{v-1}))) \end{bmatrix}'$$

and

$$\mathbb{E}(I_v I_v') = \begin{bmatrix} \mathbb{E}((y_{i,t_i} - \mathbb{P}(y_{i,t_i} | \Omega_{v-1}))(y_{k,t_k} - \mathbb{P}(y_{k,t_k} | \Omega_{v-1}))) \end{bmatrix}_{\{i=1, \dots, K; k=1, \dots, K\}}.$$

In order to obtain (12) we need to calculate $\mathbb{E}(y_{j,t_j} (y_{i,t_i} - \mathbb{P}(y_{i,t_i} | \Omega_{v-1})))$ and $\mathbb{E}((y_{i,t_i} - \mathbb{P}(y_{i,t_i} | \Omega_{v-1}))(y_{k,t_k} - \mathbb{P}(y_{k,t_k} | \Omega_{v-1})))$.

Given the model (4) we can write

$$\begin{aligned} y_{j,t_j} &= \lambda_j f_{t_j} + \xi_{j,t_j}, \\ y_{j,t_j} - y_{j,t_j | \Omega_{v-1}} &= \lambda_j (f_{t_j} - f_{t_j | \Omega_{v-1}}) + \xi_{j,t_j}. \end{aligned}$$

Therefore we have

$$\begin{aligned}
\mathbb{E}\left(y_{j,t_j} (y_{i,t_i} - y_{i,t_i|\Omega_{v-1}})\right) &= \lambda_j \mathbb{E}\left(f_{t_j} (f_{t_i} - f_{t_i|\Omega_{v-1}})'\right) \lambda'_i + \mathbb{E}\left(\xi_{j,t_j} (f_{t_i} - f_{t_i|\Omega_{v-1}})'\right) \lambda'_i \\
&= \lambda_j \mathbb{E}\left(\left(f_{t_j} - f_{t_j|\Omega_{v-1}}\right) (f_{t_i} - f_{t_i|\Omega_{v-1}})'\right) \lambda'_i + \lambda_j \mathbb{E}\left(f_{t_j|\Omega_{v-1}} (f_{t_i} - f_{t_i|\Omega_{v-1}})'\right) \lambda'_i \\
&= \lambda_j \mathbb{E}\left(\left(f_{t_j} - f_{t_j|\Omega_{v-1}}\right) (f_{t_i} - f_{t_i|\Omega_{v-1}})'\right) \lambda'_i
\end{aligned}$$

and

$$\mathbb{E}\left((y_{i,t_i} - y_{i,t_i|\Omega_{v-1}}) (y_{k,t_k} - y_{k,t_k|\Omega_{v-1}})'\right) = \lambda_i \mathbb{E}\left((f_{t_i} - f_{t_i|\Omega_{v-1}}) (f_{t_k} - f_{t_k|\Omega_{v-1}})'\right) \lambda'_k + \mathbb{E}\left(\xi_{i,t_i} \xi'_{k,t_k}\right).$$

In the case that $t_j = t_i$ the expectation $\mathbb{E}\left(\left(f_{t_j} - f_{t_j|\Omega_{v-1}}\right) (f_{t_i} - f_{t_i|\Omega_{v-1}})'\right)$ is returned by the Kalman smoother. To obtain the expectations for $t_j \neq t_i$ one can augment the vector of states by appropriate number of their lags.