

# Estimation of the Monthly Unemployment Rate through Structural Time Series Modelling in a Rotating Panel Design

*Jan van den Brakel and Sabine Krieg*

The views expressed in this paper are those of the author(s)  
and do not necessarily reflect the policies of Statistics Netherlands

**Discussionpaper (08003)**



## Explanation of symbols

.	= data not available
*	= provisional figure
x	= publication prohibited (confidential figure)
–	= nil or less than half of unit concerned
–	= (between two figures) inclusive
0 (0,0)	= less than half of unit concerned
blank	= not applicable
2005-2006	= 2005 to 2006 inclusive
2005/2006	= average of 2005 up to and including 2006
2005/'06	= crop year, financial year, school year etc. beginning in 2005 and ending in 2006
2003/'04–2005/'06	= crop year, financial year, etc. 2003/'04 to 2005/'06 inclusive

Due to rounding, some totals may not correspond with the sum of the separate figures.

*Publisher*  
Statistics Netherlands  
Prinses Beatrixlaan 428  
2273 XZ Voorburg

**second half of 2008:**  
Henri Faasdreef 312  
2492 JP The Hague

*Prepress*  
Statistics Netherlands - Facility Services

*Cover*  
TelDesign, Rotterdam

*Information*  
Telephone .. +31 88 570 70 70  
Telefax .. +31 70 337 59 94  
Via contact form: [www.cbs.nl/information](http://www.cbs.nl/information)

*Where to order*  
E-mail: [verkoop@cbs.nl](mailto:verkoop@cbs.nl)  
Telefax .. +31 45 570 62 68

*Internet*  
<http://www.cbs.nl>

ISSN: 1572-0314

### *Summary:*

*In this paper a multivariate structural time series model is described that accounts for the panel design of the Dutch Labour Force Survey and is applied to estimate monthly unemployment rates. Compared to the generalized regression estimator, this approach results in a substantial increase of the accuracy due to a reduction of the standard error and the explicit modelling of the bias between the subsequent waves.*

*Keywords: Small Area Estimation, Rotation Group Bias, Survey Errors*

## **1. Introduction**

The Dutch Labour Force Survey (LFS) is based on a rotating panel design. Each month a sample of addresses is drawn and data are collected by means of computer assisted personal interviewing of the residing households. The sampled households are re-interviewed by telephone four times at quarterly intervals. The estimation procedure of this survey is based on the generalized regression (GREG) estimator, developed by Särndal e.a. (1992).

GREG estimators are widely applied by national statistical institutes. Due to the following properties, GREG estimators are very attractive to produce official releases in a regular production environment. First, GREG estimators are approximately design-unbiased, which provides a form of robustness in the case of large sample sizes. These estimators are derived from a linear regression model that specifies the relationship between the values of a certain target parameter and a set of auxiliary variables for which the totals in the finite target population are known. If this linear regression model explains the variation of the target variable reasonably well, then this might reduce the design variance as well as the bias due to selective nonresponse, Särndal and Swenson (1987), Bethlehem (1988), and Särndal and Lundström (2005). Model misspecification, on the other hand, might result in an increase of the design variance but the point estimates remain approximately design unbiased. Second, GREG estimators are often used to produce one set of weights for the estimation of all target parameters of a multi-purpose sample survey. This is not only convenient but also enforces consistency between the marginal totals of different publication tables.

There are two major problems with the rotating panel design of the LFS and the way that the GREG estimator is applied in the estimation procedure. First there are substantial systematic differences between the subsequent waves of the panel due to mode- and panel effects. This is a well-known problem for rotating panel designs, and is in the literature referred to as rotation group bias (RGB), see Bailer (1975). In the LFS, the trend of the unemployment rate in the subsequent waves is substantially

smaller compared to the first wave. There are also systematic differences between the seasonal effects of the subsequent waves.

A second problem is that the monthly sample size of the LFS is too small to rely on the GREG estimator to produce official statistics about the monthly employment and unemployment. GREG estimators have relatively large design variance in the case of small sample sizes. Therefore, in the LFS, each month the samples observed in the preceding three months are used to estimate moving averages about the labour market situation. The major drawback of the use of a moving average is that the real monthly seasonal pattern in the unemployment rate is smoothed out and thus biased. Also structural changes in the unemployment appear delayed in the series of the published figures.

Since the monthly sample size is too small to apply design-based or direct survey estimators, model-based estimation procedures might be used to produce sufficiently reliable statistics. In the case of continuously conducted surveys, a structural time series model can be applied to use information from preceding samples to improve the accuracy of the estimates. This model can be extended to account for the RGB and the autocorrelation between the different panels of the LFS. This approach makes efficient use of the rotating panel design of the LFS in estimating monthly figures about the labour market, and is originally proposed by Pfeffermann (1991) and Pfeffermann e.a. (1998). These techniques are applied in this paper to estimate the monthly unemployment rate of the LFS. Other references to authors that apply time series models to develop estimates for periodic surveys are Scott and Smith (1974), Scott e.a. (1977), Tam (1987), Binder and Dick (1989), Bell and Hillmer (1990), Tiller (1992), Rao and Yu (1994), Pfeffermann and Burck (1990), Pfeffermann and Bleuer (1993), and Pfeffermann and Tiller (2006).

Composite estimators can be considered as an alternative. They are developed under the traditional design-based approach, to use information observed in previous periods from periodic surveys with a rotating panel design, to improve the precision of level and change estimates. Some key references to composite estimators are Hansen e.a. (1953), Rao and Graham (1964), Gurney and Daly (1965), Singh (1996), Gambino e.a. (2001), Singh e.a. (2001) and Fuller and Rao (2001). The use of composite estimators for the LFS is studied in a separate project, see Boonstra (2007).

In section 2, the survey design of the LFS is summarised. A structural time series model that accounts for the rotating panel design of the LFS is described in sections 3 and 4. The results are detailed in section 5. Some general remarks are made in section 6.

## **2. The Dutch Labour Force Survey**

### **2.1 Sample design**

The objective of the Dutch LFS is to provide reliable information about the labour market. Each month a sample of addresses is selected from which during the data collection households are identified that can be regarded as the ultimate sampling units. The target population of the LFS consists of the non-institutionalised population aged 15 years and over residing in the Netherlands. The sampling frame is a list of all known occupied addresses in the Netherlands, which is derived from the municipal basic registration of population data. The LFS is based on a stratified two-stage cluster design of addresses. Strata are formed by geographical regions. Municipalities are considered as primary sampling units and addresses as secondary sampling units. All households residing at an address, up to a maximum of three, are included in the sample (in the Netherlands, there is generally one household per address). Since most target parameters of the LFS concern people aged 15 through 64 years, addresses with only persons aged 65 years and over are undersampled.

In October 1999, the LFS changed from a continuous survey to a rotating panel design. In the first wave, data are collected by means of computer assisted personal interviewing (CAPI). Interviewers are working on the data collection of the LFS in areas around their place of residence. For all members of the selected households, demographic variables are observed. For the target variables only persons aged 15 years and over are interviewed. When a household member cannot be contacted, proxy interviewing is allowed by members of the same household. Households in which one or more of the selected persons do not respond for themselves or in a proxy interview, are treated as nonresponding households. The respondents aged 15 through 64 years are re-interviewed four times at quarterly intervals. In these four subsequent waves, data are collected by means of computer assisted telephone interviewing (CATI). During these re-interviews a condensed questionnaire is applied to establish changes in the labour market position of the respondents. Proxy interviewing is also allowed during these re-interviews. Commencing the moment that the LFS is conducted as a rotating panel design, the gross sample size has averaged about 8.000 addresses monthly, with about 65% completely responding households.

### **2.2 Rotation group bias**

The rotating panel design, described in section 2.1, results in systematic differences between the estimates of the unemployment rate of the successive waves in one time period. In the literature, this phenomenon is known as RGB, see e.g. Bailar (1975), Kumar and Lee (1983) and Pfeffermann (1991). The RGB in the LFS results in a systematic underestimation of the level of the unemployment rate in the CATI

waves but also in systematic differences between the seasonal patterns. The RGB is a consequence of the following strongly confounded factors:

- Changes in the population between the subsequent waves, due to e.g. migration and immigration.
- Selective nonresponse between the subsequent waves, i.e. panel attrition.
- Systematic differences between the populations that are reached with the CAPI and CATI modes. It is anticipated that these differences are relatively small, since the telephone number is asked during the first interview. As a result, secret numbers and cell-phone numbers are also called.
- Mode-effects, i.e. systematic differences in the data due to the fact that the interviews are conducted by telephone instead of face to face. Under the CAPI mode the interview speed is lower, respondents are more engaged with the interview and are more likely to exert the required cognitive effort to answer questions carefully. Also less socially desirable answers are obtained under the CAPI mode due to the personal contact with the interviewer. As a result, less measurement errors are expected under the CAPI mode (Holbrook e.a., 2003, and Roberts, 2007). Van den Brakel (2007) describes an experiment where the CAPI and CATI data collection modes are compared in the first wave of the LFS. It follows that the estimated unemployment rate is significantly smaller under the CATI mode.
- The fraction of proxy interviews is larger under the CATI mode (Van den Brakel, 2007). This might result in an increased amount of measurement errors.
- Effects due to differences between the CAPI questionnaire and the CATI questionnaire. The CATI questionnaire is a strongly condensed version of the CAPI questionnaire since the re-interviews focus on changes in the labour market position of the respondents.
- Panel effects, i.e. systematic changes in the behaviour of the respondents in the panel. For example, questions about activities to find a job in the first wave might increase the search activities of the unemployed respondents in the panel. Respondents might also adjust their answers in the subsequent waves systematically, since they learn how to keep the routing through the questionnaire as short as possible.

It is assumed that the estimates based on the first wave are the most reliable, since CAPI generally results in a higher data quality and the first wave does not suffer from the panel effects mentioned above. In order to minimize the effects of the RGB, the second, third, fourth and fifth waves are currently calibrated to the first wave as will be described in section 2.3.

### 2.3 Regular estimation procedure

The weighting procedure of the LFS is based on the GREG estimator (Särndal e.a., 1992). The inclusion probabilities reflect the sampling design described above as well as the different response rates between geographical regions. The weighting scheme is based on a combination of different socio-demographic categorical variables. The integrated method for weighting persons and families of Lemaître and Dufour (1987) is applied to obtain equal weights for persons belonging to the same household. Finally, a bounding algorithm proposed by Huang and Fuller (1978) is applied to avoid negative weights.

Target parameters about the employment and unemployment are defined as population totals or as ratios of two population totals. The unemployment rate, which is investigated in this paper, is defined as the ratio of the total unemployment and the total labour force. This population parameter is estimated as the ratio of the GREG estimates for the total unemployed labour force and the total labour force.

The monthly sample size of the LFS is too small to publish reliable monthly figures using the GREG estimator. Therefore each month estimates about the employment and unemployment for the preceding three months are published. The rotation scheme, described in section 2.1, implies that the sample data obtained with the five waves in three successive months are based on unique households obtained from 15 independent samples.

In an attempt to correct for the RGB, the following weighting procedure is used in the regular estimation procedure. The most important steps are summarized here. First, the inclusion weights of each CATI wave of the sample data from the preceding three months are calibrated with the GREG estimator according to the labour force status observed in the first wave using the following scheme:

$$LabourForceStatus \times Age(3).$$

*LabourForceStatus* is a classification in ten classes. *Age(i)* is a classification in *i* classes. In the next step, the calibrated weights of the four CATI waves and the inclusion weights of the CAPI wave are calibrated with the GREG estimator, using the following weighting scheme:

$$\begin{aligned} & Gender \times Ethnicity(8) \times Age(2) + Gender \times Age(21) + \\ & Age(5) \times Marital Status + Region(44) \times Age(5)Gender + \\ & Region(82) \times Age(3)Gender + Age(7)GenderCAPICATI. \end{aligned} \quad (2.1)$$

*Age(5)Gender* is a classification in eight classes which is based on five age classes and the second, third and fourth age class is itemized to gender. *Age(3)Gender* is a classification in four classes which is based on three age classes and the second age class is itemized to gender. *Region(i)* is a classification in *i* geographical regions. *Marital Status* is a classification in two classes which distinguish between married and not married. *Ethnicity(8)* is a classification in eight classes. For *Age(7)GenderCAPICATI*, it is assured that the ratio between CAPI and CATI is 3:7

for seven age classes that are itemized to gender with the exception of the first age class. This estimation procedure is conducted with the software package *Bascula* (Nieuwenbroek and Boonstra, 2002).

Since this weighting procedure hardly corrects for the RGB, an additional rigid correction is applied. For the most important parameters the ratio between the estimates based on CAPI only and the estimates based on all waves is computed using the data of 12 preceding quarters. Estimates for the preceding three months are multiplied by this ratio to correct for RGB. More information about this procedure can be found in Cuppen and Martinus (2001) and CBS documentation about the Dutch LFS.

#### 2.4 Monthly GREG estimates based on monthly data

In section 3, a structural time series model is developed to estimate the monthly unemployment rate. The input data for this time series model are the GREG estimates for the monthly unemployment rate using the monthly sample data of the separate waves. Let  $\theta_t$  denote the true but unknown unemployment rate for month  $t$ . Now  $Y_t^{t-j}$  denotes the GREG estimate of the unemployment rate of month  $t$ , based on the sample which entered the panel in month  $t-j$ . For the period of January 2001 until December 2006 each month five independent GREG estimates for the same parameter  $\theta_t$  are produced, using the five separate waves that are observed each month, i.e.  $Y_t^{t-j}$  for  $j = 0, 3, 6, 9, 12$ . The unemployment rate is estimated as

$$Y_t^{t-j} = \frac{t_{y,t}^{t-j}}{t_{z,t}^{t-j}},$$

with  $t_{y,t}^{t-j}$  and  $t_{z,t}^{t-j}$  the GREG estimates for the unemployed labour force and the labour force at time  $t$ , based on the sample that entered the panel at  $t-j$ .

The sample size that is available for the separate monthly waves, requires that weighting scheme (2.1) for the GREG estimator is reduced. For the GREG estimates of the CAPI waves, i.e.  $Y_t^t$ , the following scheme is used

$$\begin{aligned} &Age(5)Gender + Region(44) + Gender \times Age(21) + \\ &Age(5) \times Marital Status + Ethnicity(8). \end{aligned} \tag{2.2}$$

For the GREG estimates of the four CATI waves, i.e.  $Y_t^{t-j}$  for  $j = 3, 6, 9, \text{ and } 12$ , the following scheme is used:

$$\begin{aligned} &Age(5)Gender + Region(44) + Gender + Age(21) + \\ &Marital Status + Ethnicity(3). \end{aligned} \tag{2.3}$$

*Ethnicity(3)* is a classification in three classes. The estimates based on the CATI data are not adjusted to correct for RGB, since a multivariate time series model is applied to correct for this bias.



The estimates for the monthly unemployment rate obtained with the structural time series approach will be compared with monthly estimates based on the GREG estimator using the data observed in the five waves. For this comparison a slightly simplified version of the procedure described in subsection 2.3 is applied to combine the data observed in the different waves to obtain monthly GREG estimates. First, a weighted mean of all GREG estimates of the unemployment rate of month  $t$  is computed as

$$Y_t^{CAPTI} = \frac{3}{10} Y_t^t + \frac{7}{10} \left( \frac{n_t^{t-3} Y_t^{t-3} + n_t^{t-6} Y_t^{t-6} + n_t^{t-9} Y_t^{t-9} + n_t^{t-12} Y_t^{t-12}}{n_t^{t-3} + n_t^{t-6} + n_t^{t-9} + n_t^{t-12}} \right) \quad (2.4)$$

where  $n_t^{t-j}$  is the net sample size in month  $t$  of the sample which entered the panel in month  $t-j$ . Then a correction factor based on the preceding three years is computed as:

$$c_t = \frac{\sum_{j=0}^{35} Y_{t-j}^{t-j}}{\sum_{j=0}^{35} Y_{t-j}^{CAPTI}}. \quad (2.5)$$

Finally, the corrected CAPTI-estimate is computed:

$$Y_t^{CAPTI,c} = c_t Y_t^{CAPTI}. \quad (2.6)$$

Because the series start at January 2001,  $c_t$  can be computed from December 2003. To get a corrected GREG estimate for all months,  $c_{December2003}$  is used in formula (2.6) for the periods preceding December 2003.

For the computation of the standard error of (2.6), the variance of the correction factor is neglected. The variance for (2.6) is approximated by

$$\text{var}(Y_t^{CAPTI,c}) = c_t^2 \left[ \left( \frac{3}{10} \right)^2 \text{var}(Y_t^t) + \left( \frac{7}{10} \right)^2 A \right], \quad (2.7)$$

with

$$A = \frac{(n_t^{t-3})^2 \text{var}(Y_t^{t-3}) + (n_t^{t-6})^2 \text{var}(Y_t^{t-6}) + (n_t^{t-9})^2 \text{var}(Y_t^{t-9}) + (n_t^{t-12})^2 \text{var}(Y_t^{t-12})}{(n_t^{t-3} + n_t^{t-6} + n_t^{t-9} + n_t^{t-12})^2}$$

No correlation terms between  $Y_t^{t-j}$  and  $Y_t^{t-j'}$  arise in (2.7), since the estimates  $Y_t^{t-j}$ ,  $j = 0, 3, 6, 9, 12$ , are based on five independent samples. Since  $Y_t^{t-j}$  is the ratio of two GREG estimators, an expression for the estimator of the variance of  $Y_t^{t-j}$  is given by

$$\text{var}(Y_t^{t-j}) = \frac{1}{(t_{z,t}^{t-j})^2} \sum_{h=1}^H \frac{n_{h,t}^{t-j}}{n_{h,t}^{t-j} - 1} \left( \sum_{k=1}^{n_{h,t}^{t-j}} (w_k e_{k,t}^{t-j})^2 - \frac{1}{n_{h,t}^{t-j}} \left( \sum_{k=1}^{n_{h,t}^{t-j}} w_k e_{k,t}^{t-j} \right)^2 \right), \quad (2.8)$$

$$\text{with } e_{k,t}^{t-j} = \sum_{l=1}^{m_k} e_{kl,t}^{t-j},$$

$$e_{kl,t}^{t-j} = (y_{kl,t}^{t-j} - \mathbf{x}_{kl}^T \mathbf{b}_y) - Y_t^{t-j} (z_{kl,t}^{t-j} - \mathbf{x}_{kl}^T \mathbf{b}_z).$$

Here  $y_{kl,t}^{t-j}$  is a binary variable taking value one if the  $l$ -th person belonging to the  $k$ -th household that entered the sample at time  $t-j$  belongs to the unemployed labour force at time  $t$  and zero otherwise,  $z_{kl,t}^{t-j}$  a binary variable taking value one if the  $l$ -th person of the  $k$ -th household belongs to the labour force at time  $t$  and zero otherwise,  $\mathbf{x}_{kl}$  a vector with the auxiliary information of the  $l$ -th person belonging to the  $k$ -th household used in the weighting scheme of the GREG estimator,  $\mathbf{b}_y$  and  $\mathbf{b}_z$  the regression coefficient of the regression function of  $y_{kl,t}^{t-j}$  respectively  $z_{kl,t}^{t-j}$  on  $\mathbf{x}_{kl}$ ,  $w_k$  the regression weight of household  $k$ ,  $n_{h,t}^{t-j}$  the number of completely responding households of stratum  $h=1, \dots, H$  at time  $t$  of the sample that entered the panel at  $t-j$ , and  $m_k$  the number of persons aged 15 years and over belonging to the  $k$ -th household. Recall from section 2.3 that persons belonging to the same household have equal weights due to the application of the integrated method for weighting persons and families of Lemaître and Dufour (1987). Formula (2.8) is the variance estimation procedure implemented in Bascula to approximate the variance of the ratio of two GREG estimators, Nieuwenbroek and Boonstra (2002).

### 3. Time series model

Direct estimators, like the Horvitz-Thompson estimator or the GREG estimator, assume that the monthly unemployment rate  $\theta_t$  is a fixed but unknown population parameter. Under this design-based approach, an estimator for  $\theta_t$  for cross-sectional surveys only uses the data observed at time  $t$ . Data from the past are only used in the case of partially overlapping samples in a panel design, but not in the case of repeatedly conducted cross-sectional designs. Scott and Smith (1974) proposed to consider the population parameter  $\theta_t$  as a realization of a stochastic process that can be described with a time series model. Under this assumption, data observed in preceding periods  $t-1$ ,  $t-2$ , ..., can be used to improve the estimator for  $\theta_t$ , even in the case of non-overlapping sample surveys.

Due to the applied rotation pattern of the LFS, each month five independent samples are observed to estimate the population parameter  $\theta_t$ . Recall from subsection 2.4 that  $Y_t^{t-j}$  denote the GREG estimator for  $\theta_t$  based on the panel observed at time  $t$ , which entered the survey for the first time at  $t-j$ , using weighting scheme (2.2) or (2.3). Each month a vector  $\mathbf{Y}_t = (Y_t^t \ Y_t^{t-3} \ Y_t^{t-6} \ Y_t^{t-9} \ Y_t^{t-12})^T$  is observed. According to Pfeiffermann (1991), this vector can be modelled as

$$\mathbf{Y}_t = \mathbf{1}_5 \theta_t + \boldsymbol{\lambda}_t + \boldsymbol{\gamma}_t + \mathbf{e}_t, \quad (3.1)$$

with  $\mathbf{1}_5$  a five dimensional vector with each element equal to one,  $\boldsymbol{\lambda}_t = (\lambda_t^0 \lambda_t^3 \lambda_t^6 \lambda_t^9 \lambda_t^{12})^T$  and  $\boldsymbol{\gamma}_t = (\gamma_t^0 \gamma_t^3 \gamma_t^6 \gamma_t^9 \gamma_t^{12})^T$  vectors with time dependent components that account for the RGB in the trend and the RGB in the seasonal components respectively, and  $\mathbf{e}_t = (e_t^t e_t^{t-3} e_t^{t-6} e_t^{t-9} e_t^{t-12})^T$  the corresponding survey errors for each panel estimate.

Time series models for the different components in (3.1), i.e. the population parameter  $\theta_t$ , the RGB for the trend  $\boldsymbol{\lambda}_t$ , the RGB for the seasonal patterns  $\boldsymbol{\gamma}_t$ , and the survey errors  $\mathbf{e}_t$ , are developed in subsections 3.1 through 3.3.

### 3.1 Time series model for the population parameter

With a structural time series model, the population parameter  $\theta_t$  in (3.1) can be decomposed in a trend component, a seasonal component, explanatory variables and an irregular component, i.e.:

$$\theta_t = L_t + S_t + \mathbf{x}_t^T \boldsymbol{\beta}_t + \varepsilon_t, \quad (3.2)$$

where  $L_t$  denotes a stochastic trend component,  $S_t$  a stochastic seasonal component,  $\mathbf{x}_t$  a  $K$ -dimensional vector with auxiliary information,  $\boldsymbol{\beta}_t$  a  $K$ -dimensional vector with the time dependent regression coefficients  $\beta_{k,t}$  and  $\varepsilon_t$  the irregular component. For the stochastic trend component the so-called smooth trend model is used, which is defined by the following set of equations:

$$\begin{aligned} L_t &= L_{t-1} + R_{t-1}, \\ R_t &= R_{t-1} + \eta_{R,t}, \\ E(\eta_{R,t}) &= 0, \quad \text{Cov}(\eta_{R,t}, \eta_{R,t'}) = \begin{cases} \sigma_R^2 & \text{if } t = t' \\ 0 & \text{if } t \neq t'. \end{cases} \end{aligned} \quad (3.3)$$

The parameters  $L_t$  and  $R_t$  are referred to as the trend and the slope parameter respectively. The seasonal component is modelled as

$$\begin{aligned} \sum_{l=0}^{11} S_{t-l} &= \eta_{S,t}, \\ E(\eta_{S,t}) &= 0, \quad \text{Cov}(\eta_{S,t}, \eta_{S,t'}) = \begin{cases} \sigma_S^2 & \text{if } t = t' \\ 0 & \text{if } t \neq t'. \end{cases} \end{aligned} \quad (3.4)$$

Since the trend and seasonal components are modelled as stochastic processes, the values of these model parameters are allowed to change gradually over time. Similarly, the regression coefficients are assumed to be the realization of a stochastic process, to allow them to change gradually over time. A common approach is to model the regression coefficients in (3.2) as a random walk, i.e.

$$\beta_{k,t} = \beta_{k,t-1} + \eta_{\beta,k,t},$$

$$E(\eta_{\beta,k,t}) = 0, \quad Cov(\eta_{\beta,k,t}, \eta_{\beta,k',t'}) = \begin{cases} \sigma_{\beta,k}^2 & \text{if } t = t' \text{ and } k = k' \\ 0 & \text{if } t \neq t' \text{ or } k \neq k'. \end{cases} \quad (3.5)$$

The irregular component  $\varepsilon_t$  contains the unexplained variation and is modelled as a white noise process:

$$E(\varepsilon_t) = 0, \quad Cov(\varepsilon_t, \varepsilon_{t'}) = \begin{cases} \sigma_{\varepsilon}^2 & \text{if } t = t' \\ 0 & \text{if } t \neq t'. \end{cases} \quad (3.6)$$

### 3.2 Time series model for rotation group bias

The systematic differences between the trend and the seasonal components of the subsequent waves are modelled in (3.1) with  $\lambda_t$  and  $\gamma_t$ . Additional restrictions for the elements of both vectors are required to identify model (3.1). Here it is assumed that the most accurate estimate for  $\theta_t$  is obtained with the first wave, which is observed by CAPI, i.e.  $Y_t^1$ . This implies that the first component of  $\lambda_t$  and  $\gamma_t$  equals zero. Now  $\lambda_t$  measures the time dependent differences in the trend with respect to the first wave. The components of  $\lambda_t$  are defined as:

$$\lambda_t^0 = 0, \quad \lambda_t^j = \lambda_{t-1}^j + \eta_{\lambda,j,t}, \quad j = 3, 6, 9, 12, \quad (3.7)$$

$$E(\eta_{\lambda,j,t}) = 0, \quad Cov(\eta_{\lambda,j,t}, \eta_{\lambda,j',t'}) = \begin{cases} \sigma_{\lambda}^2 & \text{if } t = t' \text{ and } j = j' \\ 0 & \text{if } t \neq t' \text{ or } j \neq j'. \end{cases}$$

Furthermore  $\gamma_t$  measures the systematic differences in the seasonal components with respect to the first wave. The components of  $\gamma_t$  are defined as

$$\gamma_t^0 = 0, \quad \sum_{l=0}^{11} \gamma_{t-l}^j = \eta_{\gamma,j,t}, \quad j = 3, 6, 9, 12, \quad (3.8)$$

$$E(\eta_{\gamma,j,t}) = 0, \quad Cov(\eta_{\gamma,j,t}, \eta_{\gamma,j',t'}) = \begin{cases} \sigma_{\gamma}^2 & \text{if } t = t' \text{ and } j = j' \\ 0 & \text{if } t \neq t' \text{ or } j \neq j'. \end{cases}$$

The variance components of the random walks in (3.7) and the seasonal components in (3.8) are assumed to be equal for all waves.

### 3.3 Time series model for the survey errors

Finally a time series model for the survey errors in (3.1) is developed. From (3.1) it follows that the survey errors for the first wave are defined as  $e_t^1 = Y_t^1 - \theta_t$ . For the second, third, fourth and fifth wave, they are defined as  $e_t^{t-j} = Y_t^{t-j} - \theta_t - \lambda_t^j - \gamma_t^j$ , for  $j = 3, 6, 9, 12$ . A consequence of the rotating panel design is that the survey errors in the subsequent time periods are correlated. To account for this autocorrelation, the dependency between the survey errors of the panel observed at the last time and previous occasions is modelled with the following autoregressive relationship:

$$\begin{aligned}
e_t^t &= \eta_{e,0,t}, \\
e_t^{t-3} &= \rho_1 e_{t-3}^{t-3} + \eta_{e,3,t}, \\
e_t^{t-6} &= \rho_1 e_{t-3}^{t-6} + \rho_2 e_{t-6}^{t-6} + \eta_{e,6,t}, \\
e_t^{t-9} &= \rho_1 e_{t-3}^{t-9} + \rho_2 e_{t-6}^{t-9} + \rho_3 e_{t-9}^{t-9} + \eta_{e,9,t}, \\
e_t^{t-12} &= \rho_1 e_{t-3}^{t-12} + \rho_2 e_{t-6}^{t-12} + \rho_3 e_{t-9}^{t-12} + \rho_4 e_{t-12}^{t-12} + \eta_{e,12,t}, \\
E(\eta_{e,j,t}) &= 0, \quad Cov(\eta_{e,j,t}, \eta_{e,j',t'}) = \begin{cases} \sigma_{e,j}^2 & \text{if } t = t' \text{ and } j = j' \\ n_t^{t-j} & \\ 0 & \text{if } t \neq t' \text{ or } j \neq j'. \end{cases}
\end{aligned} \tag{3.9}$$

Here  $n_t^{t-j}$  denotes the net number of respondents in the survey at time  $t$  that entered the panel at time  $t-j$ , and  $\rho_1, \dots, \rho_4$  the autocorrelation coefficients of the AR model.

### 3.4 Final time series model for the monthly unemployment rate

The time series model for the vector with GREG estimates  $\mathbf{Y}_t$  is obtained by inserting the different components (3.2) through (3.9) into (3.1). This model uses the five monthly GREG estimates as input data to obtain model-based estimates for the monthly unemployment rate. The component for the population parameter  $\theta_t$  in (3.2), developed in subsection 3.1 makes advantages from sample information observed in the past as well as auxiliary information to improve the precision of the estimated monthly unemployment rate. The components for the RGB (3.7) and (3.8), developed in subsection 3.2, account for the systematic differences between the five monthly GREG estimates to avoid that the estimated monthly unemployment rate is biased through the non-sampling errors which are responsible for the RGB, see section 2.2. The component for the survey errors (3.9), developed in subsection 3.3, accounts for the autocorrelation between the five GREG estimates that are based on the same sample, observed with quarterly intervals.

In summary, a time series model is developed to estimate monthly unemployment rates, that makes optimal use of the available sample information from preceding periods, auxiliary information to improve the GREG estimates for the monthly unemployment rate. Furthermore, the model accounts for the rotating panel design of the LFS. Although this approach is model-based, it accounts for the complexity of the survey design of the LFS, since the GREG estimates are used as input data.

## 4. State space representation

The time series model for the five monthly GREG estimates developed in section 3 can be analysed with the Kalman filter. Therefore it has to be expressed in the so-called state space representation, see Harvey (1989) or Durbin and Koopman (2001).

A state space model consists of a measurement equation and a transition equation. The measurement equation, which is sometimes also called the signal equation, specifies how the observations depend on a linear combination of unobserved state variables, e.g. trend, seasonal, explanatory variables, RGB and the survey errors. This equation is obtained by inserting the different components for the unknown population parameter (3.2), the RGB for the trend (3.7), the RGB for the seasonal pattern (3.8) and the survey errors (3.9) into (3.1) and can be expressed as

$$\mathbf{Y}_t = \mathbf{Z}\mathbf{a}_t + \mathbf{1}_5 \varepsilon_t. \quad (4.1)$$

Here  $\mathbf{a}_t$  denotes the state vector with unobservable state variables,  $\mathbf{Z}$  a known design matrix that specifies the linear relationship between the observations and the elements of the state vector, and  $\mathbf{1}_5$  a five dimensional vector with each element equal to one. The transition equation, which is sometimes also referred to as the system equation, specifies how the state vector evolves in time:

$$\mathbf{a}_t = \mathbf{T}\mathbf{a}_{t-1} + \boldsymbol{\eta}_t, \quad (4.2)$$

with

$$E(\boldsymbol{\eta}_t) = \mathbf{0},$$

$$\text{Cov}(\boldsymbol{\eta}_t, \boldsymbol{\eta}_{t'}) = \begin{cases} \mathbf{Q}_t & \text{if } t = t' \\ \mathbf{O} & \text{if } t \neq t'. \end{cases}$$

Here  $\mathbf{0}$  and  $\mathbf{O}$  denote a vector respectively a matrix with each element zero. The state space representation of the model proposed in section 3 is obtained with (4.1) and (4.2) by taking

$$\mathbf{Y}_t = (Y_t^t \ Y_t^{t-3} \ Y_t^{t-6} \ Y_t^{t-9} \ Y_t^{t-12})^T,$$

$$\mathbf{a}_t = (\mathbf{a}_t^\theta \ \mathbf{a}_t^\lambda \ \mathbf{a}_t^\gamma \ \mathbf{a}_t^e)^T,$$

$$\mathbf{a}_t^\theta = (L_t \ R_t \ S_t \ \dots S_{t-10} \ \beta_{1,t} \ \dots \beta_{K,t}), \quad \mathbf{a}_t^\lambda = (\lambda_t^3 \ \lambda_t^6 \ \lambda_t^9 \ \lambda_t^{12}),$$

$$\mathbf{a}_t^\gamma = (\gamma_t^3 \ \dots \gamma_{t-10}^3 \ \gamma_t^6 \ \dots \gamma_{t-10}^6 \ \gamma_t^9 \ \dots \gamma_{t-10}^9 \ \gamma_t^{12} \ \dots \gamma_{t-10}^{12}),$$

$$\mathbf{a}_t^e = (\mathbf{a}_t^{e1} \ \mathbf{a}_t^{e2} \ \mathbf{a}_t^{e3} \ \mathbf{a}_t^{e4}),$$

$$\mathbf{a}_t^{e1} = (e_t^t \ e_t^{t-3} \ e_t^{t-6} \ e_t^{t-9} \ e_t^{t-12} \ e_{t-2}^{t-2} \ e_{t-2}^{t-5} \ e_{t-2}^{t-8} \ e_{t-2}^{t-11} \ e_{t-1}^{t-1} \ e_{t-1}^{t-4} \ e_{t-1}^{t-7} \ e_{t-1}^{t-10}),$$

$$\mathbf{a}_t^{e2} = (e_{t-5}^{t-5} \ e_{t-5}^{t-8} \ e_{t-5}^{t-11} \ e_{t-4}^{t-4} \ e_{t-4}^{t-7} \ e_{t-4}^{t-10} \ e_{t-3}^{t-3} \ e_{t-3}^{t-6} \ e_{t-3}^{t-9}),$$

$$\mathbf{a}_t^{e3} = (e_{t-8}^{t-8} \ e_{t-8}^{t-11} \ e_{t-7}^{t-7} \ e_{t-7}^{t-10} \ e_{t-6}^{t-6} \ e_{t-6}^{t-9}), \quad \mathbf{a}_t^{e4} = (e_{t-11}^{t-11} \ e_{t-10}^{t-10} \ e_{t-9}^{t-9}),$$

$$\mathbf{Z}_t = (\mathbf{Z}_t^\theta \ \mathbf{Z}_t^\lambda \ \mathbf{Z}_t^\gamma \ \mathbf{Z}_t^e),$$

$$\mathbf{Z}_t^\theta = (\mathbf{1}_5 \ \mathbf{0}_5 \ \mathbf{1}_5 \ \mathbf{O}_{5 \times 10} \ \mathbf{1}_5 \otimes \mathbf{x}_t^T),$$

$$\mathbf{Z}^\lambda = \begin{pmatrix} \mathbf{0}_4^T \\ \mathbf{I}_4 \end{pmatrix}, \quad \mathbf{Z}^\gamma = \begin{pmatrix} \mathbf{0}_4^T \\ \mathbf{I}_4 \end{pmatrix} \otimes (1 \ \mathbf{0}_{10}^T), \quad \mathbf{Z}^e = (\mathbf{I}_5 \ \mathbf{O}_{5 \times 26}),$$

$$\mathbf{T} = \text{Blockdiag}(\mathbf{T}^L \ \mathbf{T}^S \ \mathbf{I}_K \ \mathbf{T}^\lambda \ \mathbf{T}^\gamma \ \mathbf{T}^e),$$

$$\mathbf{T}^L = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{T}^S = \begin{pmatrix} -\mathbf{1}_{10}^T & -1 \\ \mathbf{I}_{10} & \mathbf{0}_{10} \end{pmatrix},$$

$$\mathbf{T}^\lambda = \mathbf{I}_4, \quad \mathbf{T}^\gamma = \mathbf{I}_4 \otimes \mathbf{T}^S,$$

$$\mathbf{T}^e = \begin{pmatrix} \mathbf{T}^{e11} & \mathbf{T}^{e12} & \mathbf{T}^{e13} & \mathbf{T}^{e14} \\ \mathbf{T}^{e21} & \mathbf{T}^{e22} & \mathbf{O}_{9 \times 6} & \mathbf{O}_{9 \times 3} \\ \mathbf{O}_{6 \times 13} & \mathbf{T}^{e32} & \mathbf{T}^{e33} & \mathbf{O}_{6 \times 3} \\ \mathbf{O}_{3 \times 13} & \mathbf{O}_{3 \times 9} & \mathbf{T}^{e43} & \mathbf{T}^{e44} \end{pmatrix},$$

$$\mathbf{T}^{e11} = \begin{pmatrix} \mathbf{0}_4^T & 0 & \mathbf{0}_4^T & \mathbf{0}_4^T \\ \mathbf{O}_{4 \times 4} & \mathbf{0}_4 & \rho_1 \mathbf{I}_4 & \mathbf{O}_{4 \times 4} \\ \mathbf{O}_{4 \times 4} & \mathbf{0}_4 & \mathbf{O}_{4 \times 4} & \mathbf{I}_4 \\ \mathbf{I}_4 & \mathbf{0}_4 & \mathbf{O}_{4 \times 4} & \mathbf{O}_{4 \times 4} \end{pmatrix}, \quad \mathbf{T}^{e12} = \begin{pmatrix} \mathbf{O}_{2 \times 3} & \mathbf{O}_{2 \times 6} \\ \rho_2 \mathbf{I}_3 & \mathbf{O}_{3 \times 6} \\ \mathbf{O}_{8 \times 3} & \mathbf{O}_{8 \times 6} \end{pmatrix},$$

$$\mathbf{T}^{e13} = \begin{pmatrix} \mathbf{O}_{3 \times 2} & \mathbf{O}_{3 \times 4} \\ \rho_3 \mathbf{I}_2 & \mathbf{O}_{2 \times 4} \\ \mathbf{O}_{8 \times 2} & \mathbf{O}_{8 \times 4} \end{pmatrix}, \quad \mathbf{T}^{e14} = \begin{pmatrix} \mathbf{0}_4 & \mathbf{O}_{4 \times 2} \\ \rho_4 & \mathbf{0}_2^T \\ \mathbf{0}_8 & \mathbf{O}_{8 \times 2} \end{pmatrix},$$

$$\mathbf{T}^{e21} = \begin{pmatrix} \mathbf{O}_{3 \times 5} & \mathbf{O}_{3 \times 3} & \mathbf{O}_{3 \times 5} \\ \mathbf{O}_{3 \times 5} & \mathbf{O}_{3 \times 3} & \mathbf{O}_{3 \times 5} \\ \mathbf{O}_{3 \times 5} & \mathbf{I}_3 & \mathbf{O}_{3 \times 5} \end{pmatrix}, \quad \mathbf{T}^{e22} = \begin{pmatrix} \mathbf{O}_{6 \times 3} & \mathbf{I}_6 \\ \mathbf{O}_{3 \times 3} & \mathbf{O}_{3 \times 6} \end{pmatrix},$$

$$\mathbf{T}^{e32} = \begin{pmatrix} \mathbf{O}_{4 \times 2} & \mathbf{O}_{4 \times 7} \\ \mathbf{I}_2 & \mathbf{O}_{2 \times 7} \end{pmatrix}, \quad \mathbf{T}^{e33} = \begin{pmatrix} \mathbf{O}_{4 \times 2} & \mathbf{I}_4 \\ \mathbf{O}_{2 \times 2} & \mathbf{O}_{2 \times 4} \end{pmatrix},$$

$$\mathbf{T}^{e43} = \begin{pmatrix} \mathbf{0}_2 & \mathbf{O}_{2 \times 5} \\ 1 & \mathbf{0}_5^T \end{pmatrix}, \quad \mathbf{T}^{e44} = \begin{pmatrix} \mathbf{0}_2 & \mathbf{I}_2 \\ 0 & \mathbf{0}_2^T \end{pmatrix},$$

$$\boldsymbol{\eta}_t = (\boldsymbol{\eta}_t^\theta \ \boldsymbol{\eta}_t^\gamma \ \boldsymbol{\eta}_t^\lambda \ \boldsymbol{\eta}_t^e)^T,$$

$$\boldsymbol{\eta}_t^\theta = (0 \ \eta_{R,t} \ \eta_{S,t} \ \mathbf{0}_{10}^T \ \eta_{\beta,1,t} \ \dots \ \eta_{\beta,K,t}), \quad \boldsymbol{\eta}_t^\lambda = (\eta_{\lambda,3,t} \ \eta_{\lambda,6,t} \ \eta_{\lambda,9,t} \ \eta_{\lambda,12,t}),$$

$$\boldsymbol{\eta}_t^\gamma = (\eta_{\gamma,3,t} \ \eta_{\gamma,6,t} \ \eta_{\gamma,9,t} \ \eta_{\gamma,12,t}) \otimes (1 \ \mathbf{0}_{10}^T),$$

$$\boldsymbol{\eta}_t^e = (\eta_{e,0,t} \ \eta_{e,3,t} \ \eta_{e,6,t} \ \eta_{e,9,t} \ \eta_{e,12,t} \ \mathbf{0}_{26}^T),$$

$$\mathbf{Q}_t = \text{Blockdiag}(\mathbf{Q}^\theta \ \mathbf{Q}^\gamma \ \mathbf{Q}^\lambda \ \mathbf{Q}_t^e),$$

$$\mathbf{Q}^\theta = \text{Diag}(0 \ \sigma_R^2 \ \sigma_S^2 \ \mathbf{0}_{10}^T \ \sigma_{\beta,1}^2 \ \dots \ \sigma_{\beta,K}^2), \quad \mathbf{Q}^\lambda = \sigma_\lambda^2 \mathbf{I}_4,$$

$$\mathbf{Q}^y = \text{Diag}[\sigma_y^2 \mathbf{1}_4^T \otimes (\mathbf{1} \mathbf{0}_{10}^T)],$$

$$\mathbf{Q}_t^e = \text{Diag}\left(\frac{\sigma_{e,0}^2}{n_t^t} \frac{\sigma_{e,3}^2}{n_t^{t-3}} \frac{\sigma_{e,6}^2}{n_t^{t-6}} \frac{\sigma_{e,9}^2}{n_t^{t-9}} \frac{\sigma_{e,12}^2}{n_t^{t-12}} \mathbf{0}_{26}^T\right).$$

Here  $\mathbf{1}$  denotes a vector with each element equal to one and  $\mathbf{I}$  the identity matrix. The subscripts for  $\mathbf{0}$ ,  $\mathbf{1}$ ,  $\mathbf{O}$ , and  $\mathbf{I}$  specify the dimensions of the vectors and the matrices.

The Kalman filter assumes that the disturbances of the measurement equations at different time periods are uncorrelated. This assumption is not met if the survey errors of the panel are incorporated in the irregular terms of the measurement equation. Therefore the survey errors are incorporated as unobserved components in the state vector and the dependency between the survey errors is explicitly modelled in the transition equation.

The transitional relationship of the survey errors is explained by Van den Brakel (2005). The transitional relations for the first five entries of  $\boldsymbol{\alpha}_t^e$  follow from (3.9). The remaining elements are included to have the same elements in  $\boldsymbol{\alpha}_t^e$  and  $\boldsymbol{\alpha}_{t-1}^e$  with a time shift of 1 and to assure that the vector  $\boldsymbol{\eta}_t^e$  is independent of past state vectors. This last property is required since the Kalman filter assumes that  $\text{Cov}(\boldsymbol{\eta}_t, \boldsymbol{\eta}_{t'}) = \mathbf{O}$  for  $t \neq t'$ .

Under the assumption of normally distributed error terms, the Kalman filter can be applied to obtain optimal estimates for the state vector  $\boldsymbol{\alpha}_t$ . Estimates for state variables for period  $t$  based on the information available up to and including period  $t$  are referred to as the filtered estimates. The filtered estimates of past state vectors can be updated, if new data become available. This procedure is referred to as smoothing and result in smoothed estimates that are based on the completely observed time series. So the smoothed estimate for the state vector for period  $t$  also accounts for the information made available after time period  $t$ . In this paper, the Kalman filter estimates for the state variables are smoothed with the fixed interval smoother. See Harvey (1989) or Durbin and Koopman (2002) for technical details.

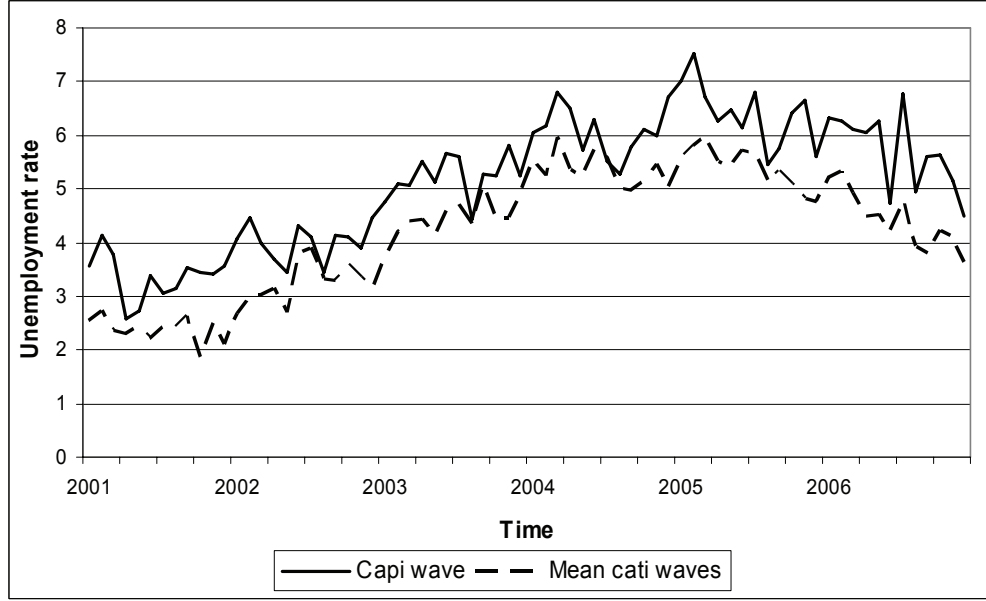
## 5. Results

### 5.1 Preliminary analyses

With the GREG estimator monthly estimates for the unemployment rate are obtained for each wave as described in section 2.4. In Figure 5.1 the unemployment rate based on the CAPI wave is compared with the average of the four CATI waves. The graph shows that the unemployment rate observed with the first wave is systematically higher than the other four waves.



Figure 5.1: RGB monthly unemployment rate based on GREG estimates



The five time series obtained with the different waves are modelled with the time series model proposed in sections 3 and 4. The analysis is conducted with software developed in Ox in combination with the subroutines of SsfPack (beta 3), see Doornik (1998) and Koopman e.a. (1999). Note that a more recent version of Ssfpack is used than the 2.2 version described in Koopman e.a. (1999). Version 3 is very appropriate for the estimation of multivariate structural time series models.

Preliminary analyses indicate that the model proposed in sections 3 and 4 can be simplified. First, the explanatory variables for the population parameter are deleted in the final model. A potential auxiliary variable is the registered unemployment. The number of persons registered as being unemployed at the Office for Employment and Income (in Dutch abbreviated as CWI) does not have a significant contribution in a model with a trend and seasonal component; see Van den Brakel and Krieg (2006). Therefore no auxiliary information is used in this paper and the state space representation of the model in section 4 can be simplified by taking

$$\begin{aligned} \mathbf{\alpha}_t^\theta &= (L_t \ R_t \ S_t \ \dots \ S_{t-10}), & \mathbf{Z}^\theta &= (\mathbf{1}_5 \ \mathbf{0}_5 \ \mathbf{1}_5 \ \mathbf{O}_{5 \times 10}), \\ \boldsymbol{\eta}_t^\theta &= (0 \ \eta_{R,t} \ \eta_{S,t} \ \mathbf{0}_{10}^T), & \mathbf{T} &= \text{Blockdiag}(\mathbf{T}^L \ \mathbf{T}^S \ \mathbf{T}^\lambda \ \mathbf{T}^\gamma \ \mathbf{T}^e), \\ \mathbf{Q}^\theta &= \text{Diag}(0 \ \sigma_R^2 \ \sigma_S^2 \ \mathbf{0}_{10}^T). \end{aligned}$$

Secondly, the estimates for the RGB of the seasonal effects in the second wave are not significantly different from zero and the RGB for the seasonal effects of the third, fourth and fifth wave are not significantly different from each other. Therefore the model is simplified by taking

$$\mathbf{\alpha}_t^\gamma = (\gamma_t \ \dots \ \gamma_{t-10})^T, \quad \mathbf{Z}^\gamma = (0 \ 0 \ 1 \ 1 \ 1)^T \otimes (\mathbf{1} \ \mathbf{0}_{10}^T),$$

$$\mathbf{T}^\gamma = \mathbf{T}^S, \quad \boldsymbol{\eta}_t^\gamma = (\eta_\gamma \mathbf{0}_{10}^T),$$

$$\mathbf{Q}^\gamma = \text{Diag}[\sigma_\gamma^2 \mathbf{0}_{10}^T].$$

Thirdly, the estimates for the AR parameters of lag two, three and four tend to zero, so the model for the survey errors (3.9) is simplified to an AR(1) model. This implies that the accompanying components in the state space model are simplified by taking

$$\boldsymbol{\alpha}_t^e = \boldsymbol{\alpha}_t^{e1}, \quad \mathbf{Z}^e = (\mathbf{I}_5 \mathbf{0}_{5 \times 8}), \quad \mathbf{T}^e = \mathbf{T}^{e11},$$

$$\boldsymbol{\eta}_t^e = (\eta_{e,0,t} \eta_{e,3,t} \eta_{e,6,t} \eta_{e,9,t} \eta_{e,12,t} \mathbf{0}_8^T),$$

$$\mathbf{Q}_t^e = \text{Diag} \left( \frac{\sigma_{e,0}^2}{n_t^t} \frac{\sigma_{e,3}^2}{n_t^{t-3}} \frac{\sigma_{e,6}^2}{n_t^{t-6}} \frac{\sigma_{e,9}^2}{n_t^{t-9}} \frac{\sigma_{e,12}^2}{n_t^{t-12}} \mathbf{0}_8^T \right).$$

## 5.2 Estimation results for the time series model

Maximum likelihood estimates for the hyperparameters, i.e. the variance components of the stochastic processes for the state variables  $\sigma_R^2, \sigma_S^2, \sigma_\varepsilon^2, \sigma_\lambda^2, \sigma_\gamma^2, \sigma_{e,0}^2, \sigma_{e,3}^2, \sigma_{e,6}^2, \sigma_{e,9}^2, \sigma_{e,12}^2$ , and the AR parameter  $\rho_1$ , are obtained using a numerical optimization procedure. The results are presented in Table 5.1.

*Table 5.1: Maximum likelihood estimates hyperparameters*

Hyperparameter	estimate
Slope (smooth trend)	0.000184
Seasonal	0.000001
RGB trend	0.000000
RGB seasonal	0.000000
Irregular component unemployment rate	0.000955
Survey error wave 1	0.341
Survey error wave 2	0.246
Survey error wave 3	0.287
Survey error wave 4	0.333
Survey error wave 5	0.278
First order auto correlation parameter for survey error	0.2

The smoothed Kalman filter estimates for the unemployment rate  $\theta_t$  are given in Figure 5.2. These are the estimates for the monthly unemployment rate, based on the smooth trend model and a seasonal component, corrected for the RGB between the five GREG estimates. The trend component is time dependent since the maximum likelihood estimate of the corresponding hyperparameter is positive (see Table 5.1). The seasonal component is almost constant since the maximum likelihood estimate of the corresponding hyperparameter is very small. The smoothed Kalman filter estimates for the trend and the seasonal component are plotted in Figures 5.3 and 5.4 respectively.

Figure 5.2: Smoothed Kalman filter estimates for the monthly unemployment rate

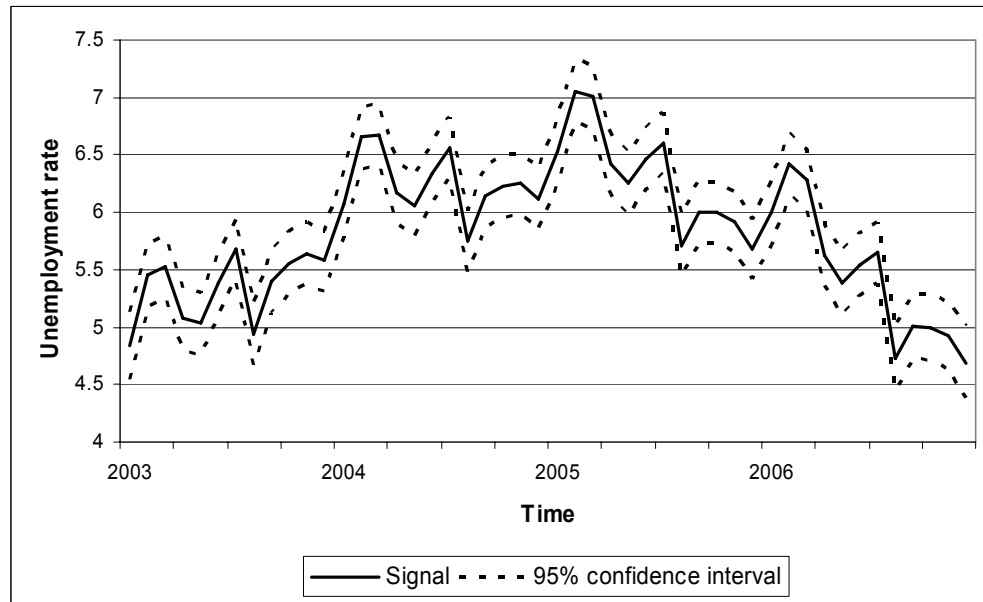


Figure 5.3: Smoothed Kalman filter estimates for the trend of the monthly unemployment rate

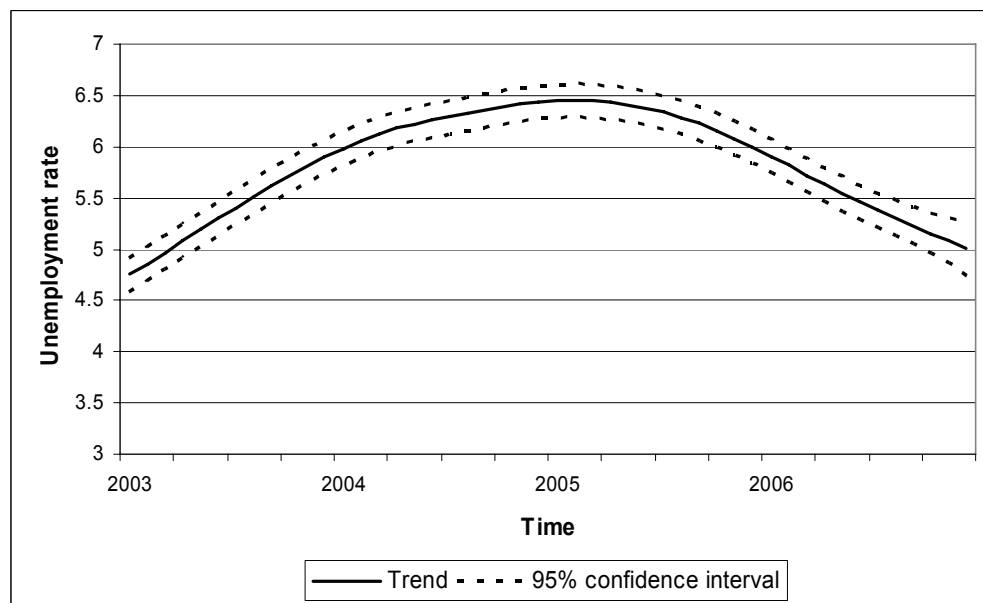
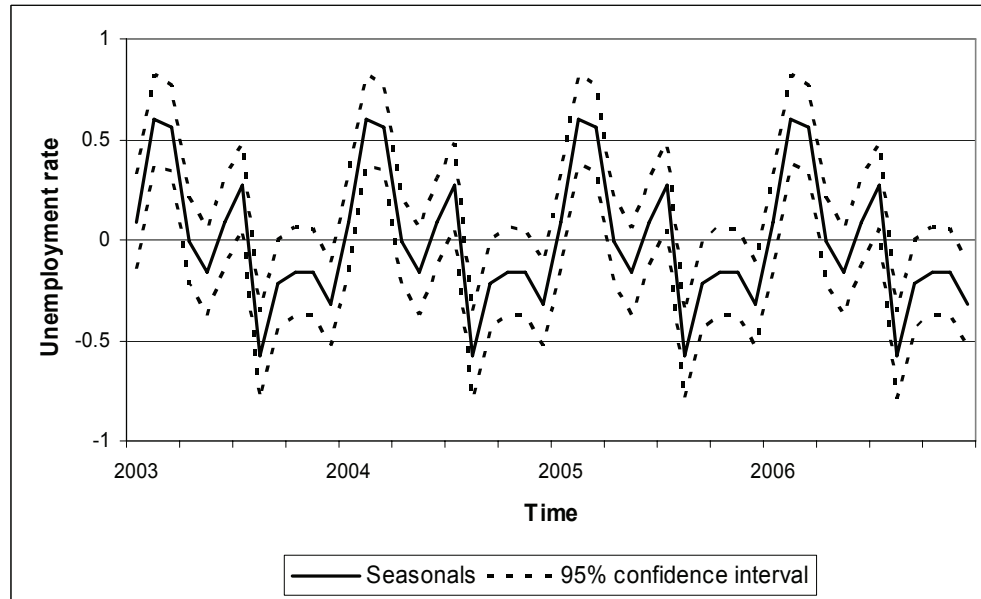


Figure 5.4: Smoothed Kalman filter estimates for the seasonal effect of the monthly unemployment rate



The Kalman filter estimates for the RGB of the trend are time independent since the maximum likelihood estimate of the corresponding hyperparameter tends to zero and is therefore put to zero. The smoothed Kalman filter estimates for the RGB are given in Table 5.2. The model beautifully detects a slightly increasing bias in the trend of the subsequent waves. The estimates for the RGB of the four CATI waves are significantly different from zero.

Table 5.2: Smoothed Kalman filter estimates RGB trend

Wave	RGB	St. error
2	-0.77	0.06
3	-0.89	0.07
4	-0.94	0.08
5	-1.11	0.07

The Kalman filter estimates for the RGB of the seasonal effects are also time independent since the maximum likelihood estimate of the corresponding hyperparameter tends to zero, which is therefore put to zero. The smoothed Kalman filter estimates are given in Figure 5.5. The smoothed Kalman filter estimates of the seasonal effects are compared with the smoothed estimates for the RGB of the seasonal effects in Figure 5.6.

Figure 5.5: Smoothed Kalman filter estimates for the RGB of the seasonal effects in the third, fourth and fifth wave

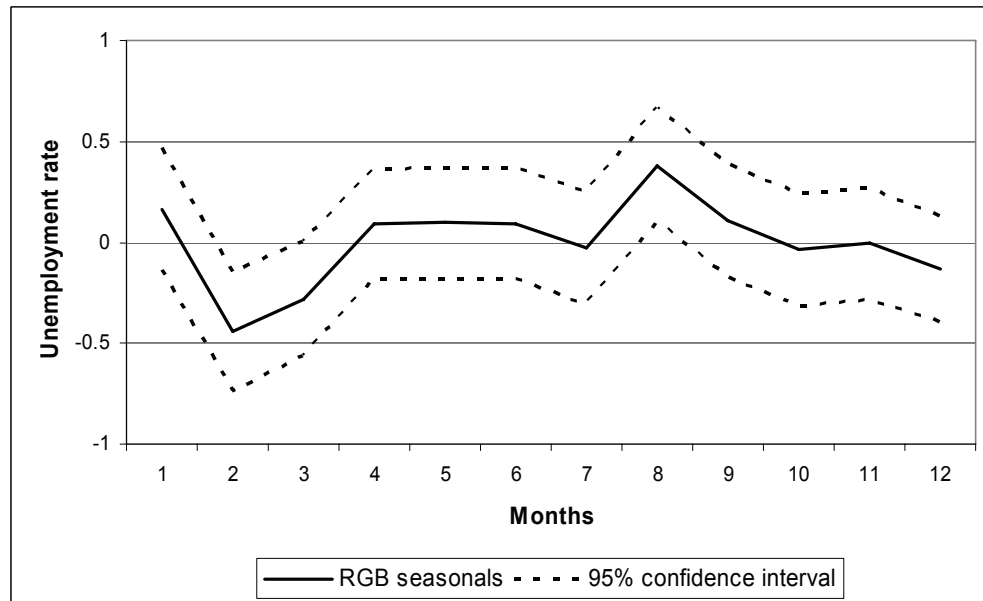
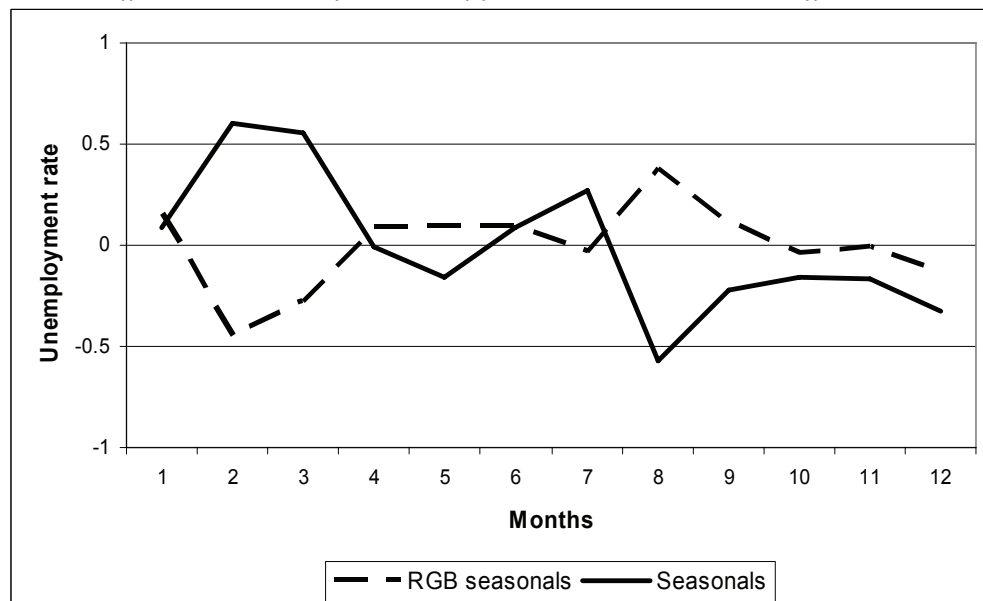


Figure 5.6: Comparison of smoothed Kalman filter estimates for the RGB of the seasonal effects in the third, fourth and fifth wave and the seasonal effects in 2006



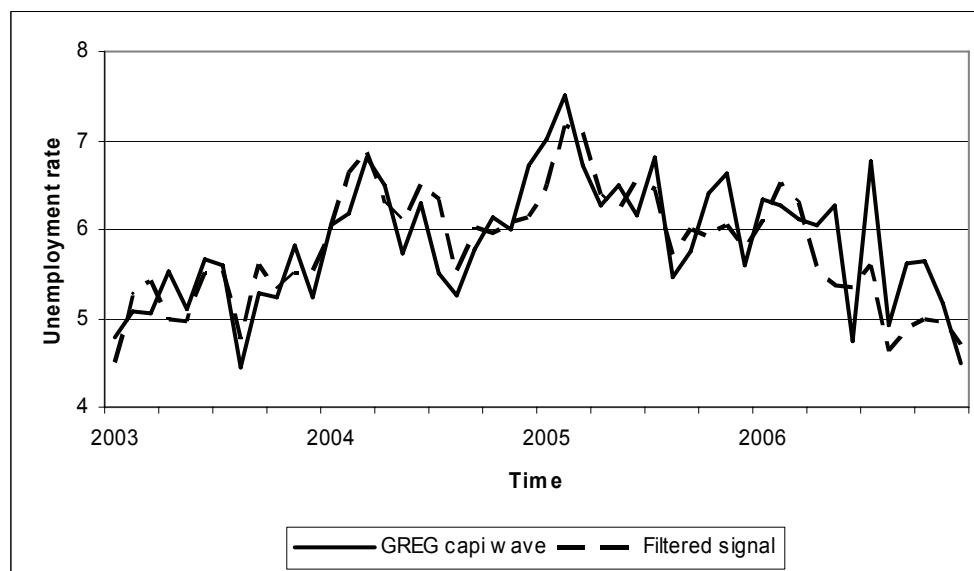
It follows from Figure 5.5 that the seasonal effects in particularly February and August in the third, fourth and fifth wave are significantly different from the first and the second wave. Figure 5.6 shows that the RGB in the seasonal effects largely nullifies the seasonal effects in February, March, and August. The seasonal effects in the last three waves are, apparently, less pronounced than in the first two waves. The different factors that contribute to the RGB in both the trend and the seasonal patterns are summarised in section 2.2.

### 5.3 Comparison with GREG estimates

In this section the monthly estimates for the unemployment rate and their standard errors obtained with the GREG estimator are compared with the filtered estimates obtained with the time series model. The filtered estimates are used since they are based on the complete set of information that would be available in the regular production process to produce a model-based estimate for the monthly unemployment rate for month  $t$ .

The GREG estimates based on the CAPI wave for the monthly unemployment rates, using weighting scheme (2.2), are compared with the filtered estimates obtained with the time series model in Figure 5.7. Some of the peaks and dips in the series of the GREG estimates are partially considered as survey errors under the structural time series model and flattened out in the filtered estimates for the series. Some of these peaks and dips are preserved since they are considered as seasonal effects under the time series model. It also follows that the filtered estimates are corrected for the RGB since the filtered series is at the same level as the GREG series for the CAPI wave. This is enforced in the way the time series model is identified, since the model parameters in (3.7) and (3.8) for the RGB for the first wave are assumed to be zero. This implies that the CATI waves are benchmarked to the outcomes of the first wave.

Figure 5.7: Filtered estimates and GREG estimates based on the CAPI wave for the monthly unemployment rate

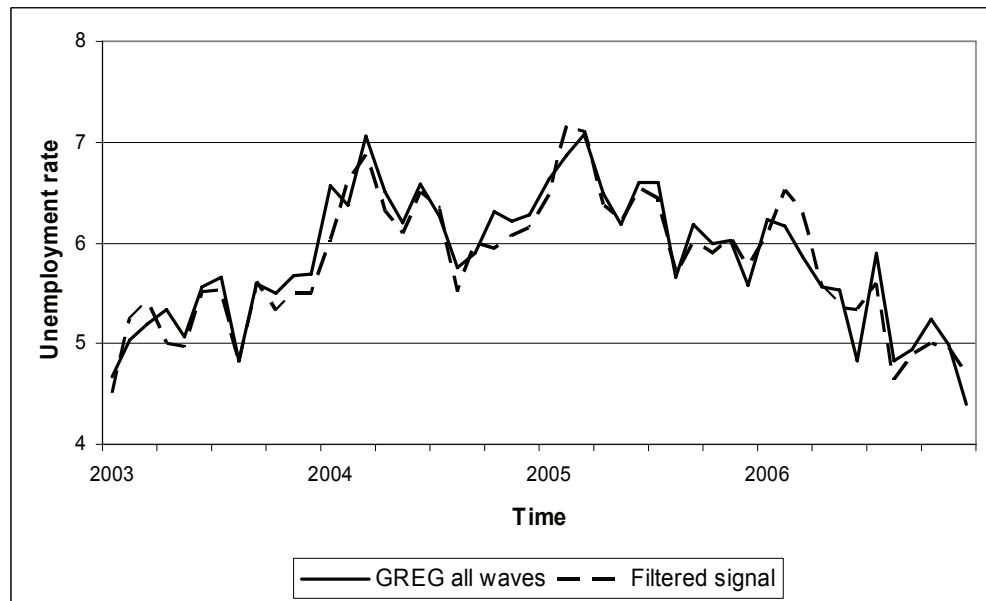


The procedure applied in the regular estimation procedure of the LFS, to combine the CATI and the CAPI waves, is also used to estimate monthly unemployment figures. This procedure is described in section 2.4. The GREG estimates for the monthly unemployment rates based on the five waves, using formula (2.6), are

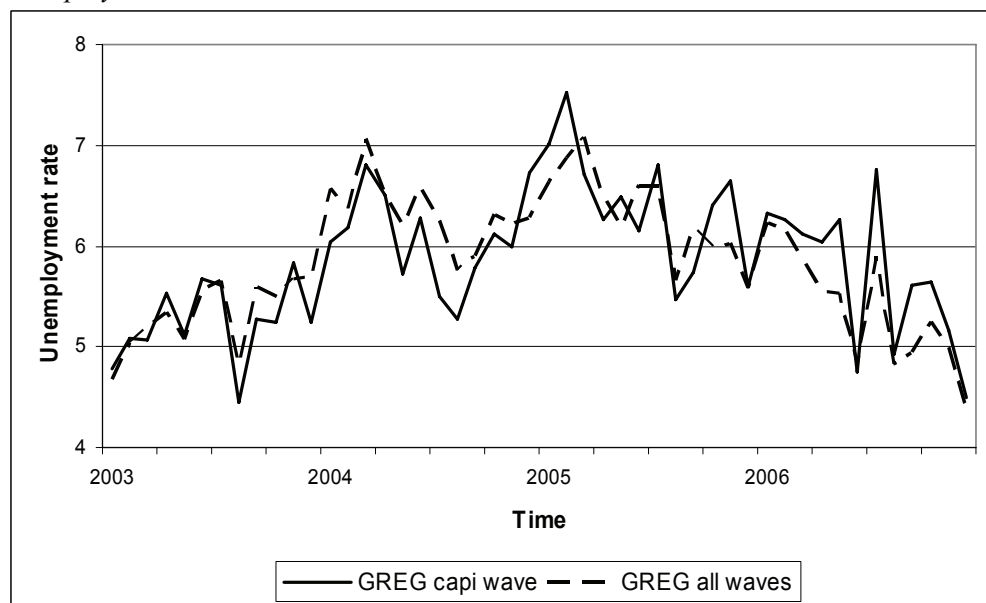
compared with the filtered estimates obtained with the time series model in Figure 5.8. Both estimates for the monthly unemployment rate follow the same level, since they are both benchmarked to the outcomes of the first wave. The GREG estimator is benchmarked in a rather rigid way using ratio (2.5), which is assumed to be constant in advance over a period of three years. The filtered estimates are benchmarked in a more subtle way through the explicit modelling of the RGB in section 3.2 through stochastic processes that are allowed to be time dependent.

The monthly GREG estimates based on all waves are also compared with the GREG estimates based on the CAPI wave in Figure 5.9.

*Figure 5.8: Filtered estimates and GREG estimates based on all waves for the monthly unemployment rate*



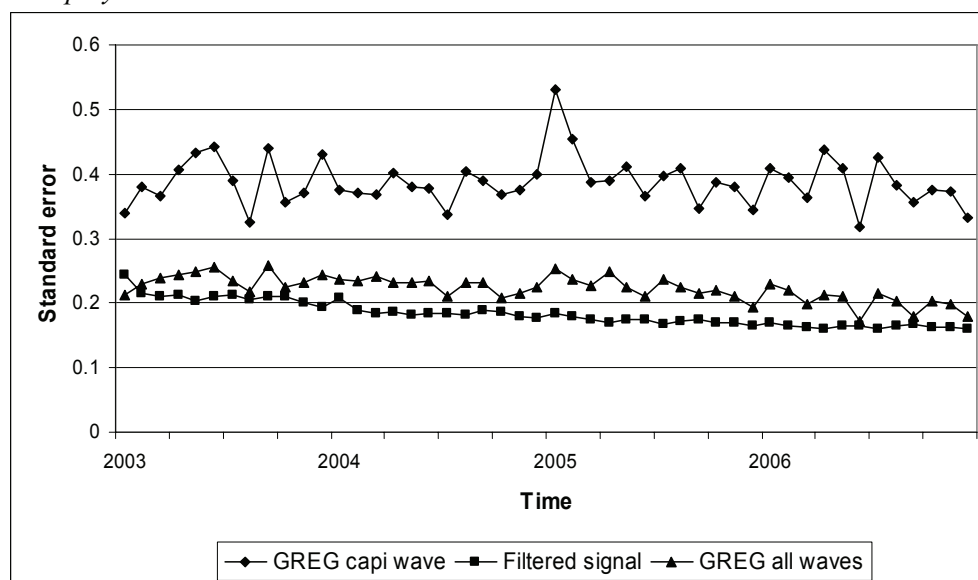
*Figure 5.9: GREG estimates based on the CAPI wave and all waves for the monthly unemployment rate*



The ratio correction applied in formula (2.6) to the GREG estimates based on all waves removes the RGB in the trend between the subsequent waves, but does not correct for the RGB in the seasonal patterns. This follows from Figure 5.8 and 5.9. The series of the GREG estimates based on all waves follows the same level as the GREG estimates based on the CAPI wave (Figure 5.9). There are, however, subtle differences between the filtered estimates obtained with the time series model and the GREG estimate based on all waves (Figures 5.8). They partially arise because some of the dips and peaks in the GREG estimates are considered as survey errors by the time series model but they are also the result of systematic differences in the seasonal patterns between the subsequent waves. The model estimates in February and March are for example larger in 2003, 2005 and 2006, and smaller in August in 2004, 2005 and 2006.

The standard errors of the monthly GREG estimates based on all waves, the CAPI wave and the filtered estimates are compared with each other in Figure 5.10. The standard errors for the GREG estimates are based on formula (2.7). The standard errors of the GREG estimates that are based on the CAPI wave are based on formula (2.8). Standard errors of the filtered estimates are obtained by the standard recursion formulas of the Kalman filter, see Harvey (1989) or Durbin and Koopman (2001).

*Figure 5.10: Standard errors of the GREG and filtered estimates for the monthly unemployment rate*



As expected, the standard errors of the GREG estimates based on all waves are smaller than the GREG estimates based on the CAPI wave, since they are based on more data. The standard errors of the filtered estimates obtained with the time series model are smaller than the GREG estimates based on all waves, since the time series model uses additional sample information from preceding periods. The standard



errors of the filtered estimates are slightly but continuously decreasing during the period 2003 to 2006. This indicates that the filter is picking up new information if new data becomes available. Smaller standard errors for the filtered estimates might be expected if more data become available.

The GREG estimates are corrected for the RGB in the trend only in a rather rigid way by applying a ratio correction which is almost constant since it is applied over a period of three years. The time series model, on the other hand, accounts for the RGB in the trend but also for the seasonal patterns in a more subtle way, since they are explicitly modelled with stochastic processes that are allowed to vary over time. This requires a rather complex model with a large amount of hyperparameters and state vectors to be estimated. The standard errors of the GREG estimator do not reflect bias that arises from the systematic differences between the seasonal patterns of the subsequent waves.

The size and complexity of the applied time series model, is large compared to the length of the series available to fit the model. The final model that is applied to a five dimensional series which is monthly observed during a period of six years contains 41 state variables. Therefore it is worthwhile to consider more parsimonious models, which might reduce the standard errors of the filtered estimates. Furthermore, the efficiency obtained by borrowing sample information from the past by relying on a time series model is illustrated more clearly if the standard error of the GREG estimates using all waves is compared with the standard error of the monthly estimates obtained with a time series model that accounts for the RGB in the trend only. Therefore a time series model without a component for the RGB in the seasonal pattern is applied to the data in an attempt to further improve the precision of the time series model estimates. This implies that  $\mathbf{a}_t^\gamma$ ,  $\mathbf{Z}^\gamma$ ,  $\mathbf{T}^\gamma$ ,  $\boldsymbol{\eta}_t^\gamma$ , and  $\mathbf{Q}^\gamma$  are deleted from the model in state space representation as described in section 4. The filtered estimates for the monthly unemployment rates based on a model with and without a component for the RGB in the seasonal pattern are compared in Figure 5.11.

The model without a component for the RGB of the seasonal effects assumes a seasonal effect for the population parameter  $\theta_t$  that is based on an average of the seasonal effects of the five waves. As a result the absolute values of the seasonal effects in February, March, and August are smaller under the simplified model, resulting in a lower estimate for the monthly unemployment rate in February and March and a larger estimate in August.

The standard errors of the filtered estimates obtained with the two time series models and the standard errors of the GREG estimates using all waves are compared in Figure 5.12. The standard error of the filtered estimates of the simplified time series model is substantially smaller than the standard error of the GREG estimates using all waves. This is the increase in precision that is obtained by using the sample

information from preceding periods through the time series model. The simplification of the time series model by ignoring the RGB for the seasonal effects, results in a reduction of the standard error at the cost of an increased bias in the seasonal effects. Under the model assumption that the estimates based on the first wave are unbiased, the time series model that accounts for the RGB in the seasonal effects is preferred, since it removes the bias in the seasonal pattern.

Figures 5.11: Filtered estimates of the monthly unemployment rate for two different time series models

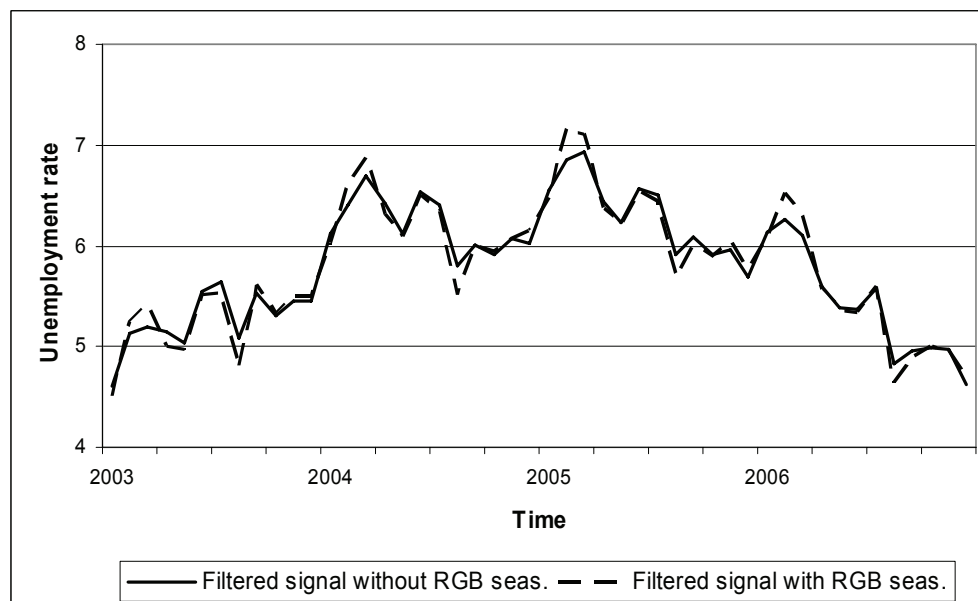
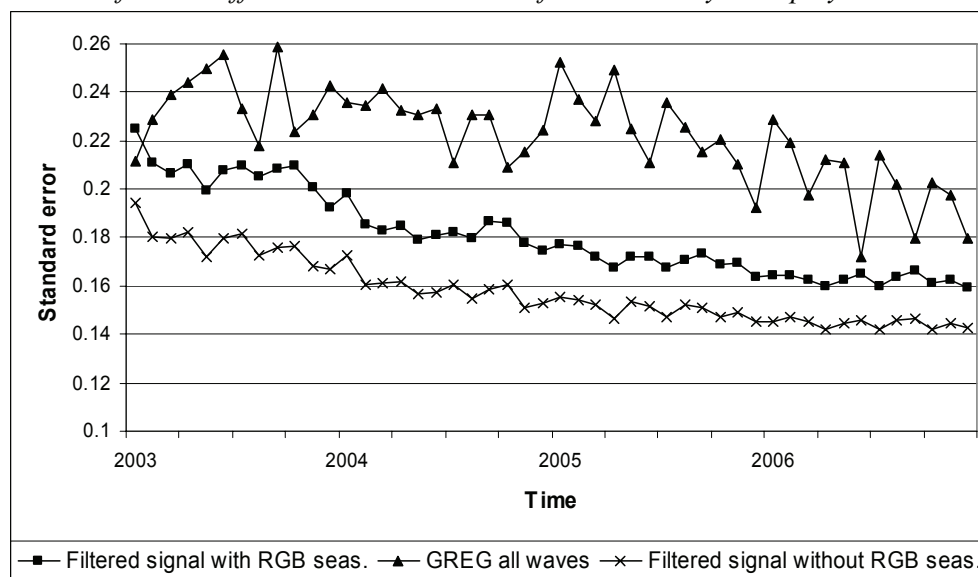


Figure 5.12: Standard errors of the GREG estimates based on all waves and filtered estimates for two different time series models for the monthly unemployment rate



## 6. Discussion

In this paper a multivariate structural time series model is applied to the monthly data of the LFS that accounts for the rotating panel design of this survey. This approach is initially proposed by Pfeffermann (1991) and extended in this paper with a component that models systematic differences in the seasonal effects between the subsequent waves. With this time series model a substantial increase of the accuracy of the monthly estimates for the unemployment rate is obtained. Firstly, the model explicitly estimates the RGB in the trend and the seasonal patterns between the first CAPI wave and the four subsequent CATI waves. As a result, estimates for the unemployment rates are corrected for this RGB. Secondly, the time series model borrows strength from data observed in preceding periods via the assumed model for the population parameter and the autocorrelation between the survey errors of the different panels.

The RGB induced by the rotating panel design is substantial. The bias in the trend results in an underestimation of the unemployment rate in the subsequent waves and its magnitude slightly decreases from -0.8 percent points in the second wave to -1.1 percent points in the fifth wave. The seasonal effect in February is about 0.5 percent points too small and in August 0.4 percent points too large in the third, fourth and fifth wave compared to the first two waves. This results in less pronounced seasonal effects in the last three waves.

The estimation procedure used in the regular survey process of the LFS is based on the GREG estimator. In this procedure the estimates for the unemployment rate are corrected with the ratio between the unemployment rate based on CAPI only and the estimates based on all waves, using the data of 12 preceding quarters. This ratio corrects for the RGB in the trend but not for the RGB in the seasonal patterns. Compared with the currently applied estimation procedure, the time series model improves the accuracy of the estimates of the unemployment rate, since it reduces the standard error and gives, under the assumption that the data obtained in the first wave are not biased, better corrections for the RGB.

A parsimonious time series model that accounts for the RGB in the trend but not for the RGB in the seasonal pattern, results in a further reduction of the standard error of the filtered estimates. This, however, results in a biased seasonal pattern in the monthly estimates of the unemployment rates. Since the standard errors of the filtered estimates obtained under this parsimonious model do not reflect this bias, a time series model that account for both the RGB in the trend and the seasonal pattern is preferred.

The time series model is identified by adopting a restriction for the RGB parameters which assumes that the first wave is observed without bias. This implies that the estimates based on the first wave are used to benchmark the subsequent waves. If

this restriction is used, then an all out effort in each part of the statistical process is required to reduce possible bias in the first wave, e.g. by using the most appropriate mode, reducing nonresponse, optimizing the weighting scheme, etc. Based on external information about the bias in the different waves, the restrictions for the RGB components might be adjusted.

The time series approach explored in this paper is appropriate to produce model-based estimates for monthly unemployment figures. Statistics Netherlands, however, is generally rather reserved in the application of model-based estimation procedures for the production of official statistics. Model misspecification might result in severely biased estimates. This bias is not reflected in the standard errors of the Kalman filter estimates. Extensive model selection and evaluation is therefore required for each separate target variable. This hampers a straightforward application of such estimation techniques, since there is generally limited time available for the analysis phase of the regular production process of official releases.

There is, on the other hand, a case for having official series that are based on model-based procedures with appropriate methodology and quality descriptions for situations where direct estimators do not result in sufficiently reliable estimates. The RGB observed under the rotating panel design of the LFS clearly illustrates the existence of non-sampling errors such as measurement errors and panel attrition. Therefore the traditional concepts that observations obtained from sampling units are true fixed values observed without error and that the respondents can be considered as a representative probability sample from the target population, generally assumed in design-based sampling theory, are not tenable under such designs. The application of direct estimators in the case of measurement errors and selective panel attrition will result in severely biased estimates. The time series model applied in this paper can be used to produce estimates that are corrected for the bias introduced by these non-sampling errors.

This estimation procedure is also applicable in situations where small sample sizes result in unacceptable large standard errors. Small sample sizes arise if official statistics are required for small domains or for short data collection periods like the monthly unemployment figures in the LFS. Most surveys conducted by national statistical institutes operate continuously in time and are based on cross-sectional or rotating panel designs. Consequently, estimation procedures based on time series models that use sample information observed in preceding periods are particularly interesting.

The time series model yields estimates for the trend and seasonal components of the population parameter. Seasonally adjusted parameter estimates and their estimation errors are therefore obtained as a by-product of this estimation procedure. Another major advantage is that this approach accounts for the autocorrelation in the survey errors due to the rotating panel design. Pfeffermann, e.a. (1998) show that ignoring

these autocorrelation, for example with the Henderson filters in X-11-ARIMA (Findley e.a. 1998), results in spurious trend estimates.

The model can be improved in several ways. Not all possible auxiliary information that is available in the register of the Office for Employment and Income has been investigated yet. Additional research about the use of registered unemployment and related variables that can be used as auxiliary information in the models is still needed. To describe the seasonal patterns, trigonometric models for the seasonal components can be considered as an alternative. Another possible improvement is detection and modelling of outliers. Furthermore the model needs to be extended to estimate monthly unemployment rates for different domains using sample information collected in the past as well as cross-sectional data from other small areas, using the approach proposed by Pfeffermann and Burck (1990) and Pfeffermann and Tiller (2006).

## References

- Bailar, B.A. (1975). The Effects of Rotation Group Bias on Estimates from Panel Surveys. *Journal of the American Statistical Association*, 70, pp. 23-30.
- Bell, W.R. and S.C. Hillmer (1990). The Time Series Approach to Estimation of Periodic Survey. *Survey Methodology*, 16, pp. 195-215.
- Bethlehem, J.G. (1988). Reduction of nonresponse bias through regression estimation. *Journal of Official Statistics*, 4, pp. 251-260.
- Binder, D.A. and J.P. Dick (1989). Modeling and Estimation for Repeated Surveys. *Survey Methodology*, 15, pp. 29-45.
- Boonstra, H.J. (2007). Samengestelde schatters voor roterende panels door ophoging. Unpublished report, BPA nr: DMK-DMH-2007-09-20-HBTA, Statistics Netherlands, Heerlen (in Dutch).
- Brakel, J.A. van den (2007). Design-based Analysis of Embedded Experiments with Applications in the Dutch Labour Force Survey. Discussion paper 07009, <http://www.cbs.nl/NR/rdonlyres/7010D33B-8C66-4028-9DB1-0865E90B39E9/0/200709x10pub.pdf>, Statistics Netherlands, Heerlen.
- Brakel, J.A. van den (2005). Small Area Estimators for the Dutch Labour Force Survey using Structural Time Series Models. Unpublished research paper, BPA nr: TMO-R&D-2005-05-02-JBRL, Statistics Netherlands, Heerlen.
- Brakel, J.A. van den, S. Krieg (2006). Kleinedomeinschatters bij de Enquête Beroepsbevolking via Tijdreeksmodellen. Unpublished research paper (in Dutch), BPA nr: TMO-R&D-2006-01-20-JBRL, Statistics Netherlands, Heerlen.

- Cuppen, M. and G.H. Martinus (2001). Weegmodel voor de Enquête Beroepsbevolking. Unpublished research paper, BPA nr: 1948-01-TMO, Statistics Netherlands, Heerlen (in Dutch).
- Doornik, J.A. (1998). *Object-Oriented Matrix Programming using Ox 2.0*. London: Timberlake Consultants Press.
- Durbin, J. and S.J. Koopman (2001). *Time series analysis by state space methods*. Oxford: Oxford University Press.
- Findley, D.F., B.C. Monsell, W.R. Bell, M.C. Otto, and B.C. Chen (1998). New capabilities and methods of the X-12-ARIMA Seasonal Adjustment Program. *Journal of Business and Economic Statistics*, 16, pp. 127-176 (with Discussion).
- Fuller, W.A. and J.N.K. Rao (2001). A Regression Composite Estimator with Application to the Canadian Labour Force Survey. *Survey Methodology*, 27, pp. 45-51.
- Gambino, J., B. Kennedy, and M.P. Singh, (2001). Regression Composite Estimation for the Canadian Labour Force Survey: Evaluation and Implementation. *Survey Methodology*, 27, pp. 65-74.
- Gurney, M., and J.F. Daly (1965). A multivariate approach to estimation in periodic sample surveys. *Proceedings of the Section on Social Statistics*, American Statistical Association, pp. 242-257.
- Hansen, M.H., W.N. Hurwitz, and W.G. Meadow, (1953). *Sample Survey Methods and Theory*, 2. New York: Wiley.
- Harvey, A.C. (1989). *Forecasting, Structural Time Series Models and the Kalman Filter*. Cambridge: Cambridge University Press.
- Holbrook, A.L., M.C. Green, and J.A. Krosnick (2003). Telephone versus Face-to-Face Interviewing of National Probability Samples with Long Questionnaires. *Public Opinion Quarterly*, 67, pp. 79-125.
- Huang, E.T. and W.A. Fuller (1978). Nonnegative Regression Estimation for Survey Data, *Proceedings of the Section on Social Statistics*, American Statistical Association 1978, pp. 300-303.
- Koopman, S.J., N. Shephard and J.A. Doornik (1999). Statistical Algorithms for Models in State Space using SsfPack 2.2. *Econometrics Journal*, 2, pp. 113-166.
- Kumar, S. and H. Lee. (1983). Evaluation of Composite Estimation for the Canadian Labour Force Survey. *Survey Methodology*, 9, pp. 178-201.
- Lemaître, G. and J. Dufour (1987). An Integrated Method for Weighting Persons and Families. *Survey Methodology*, 13, pp. 199-207.
- Nieuwenbroek, N. and H.J. Boonstra (2002). Bascula 4.0 Reference Manual, BPA nr: 279-02-TMO, Statistics Netherlands, Heerlen.

- Pfeffermann, D. (1991). Estimation and Seasonal Adjustment of Population Means Using Data from Repeated Surveys. *Journal of Business & Economic Statistics*, 9, pp. 163-175.
- Pfeffermann, D. and S.R. Bleuer (1993). Robust Joint Modelling of Labour Force Series of Small Areas. *Survey Methodology*, 19, pp. 149-163.
- Pfeffermann, D. and L. Burck (1990). Robust Small Area Estimation Combining Time Series and Cross-Sectional Data. *Survey Methodology*, 16, pp. 217-237.
- Pfeffermann, D., M. Feder and D. Signorelli (1998). Estimation of Autocorrelations of Survey Errors with Application to Trend Estimation in Small Areas. *Journal of Business & Economic Statistics*, 16, pp. 339-348.
- Pfeffermann, D. and R. Tiller (2006). Small Area Estimation with State Space Models Subject to Benchmark Constraints. *Journal of the American Statistical Association*, 101, pp. 1387-1397.
- Rao, J.N.K., and J.E. Graham (1964). Rotating Designs for Sampling on Repeated Occasions. *Journal of the American Statistical Association*, 59, pp. 492-509.
- Rao, J.N.K., and M. Yu (1994), Small Area Estimation by Combining Time Series and Cross-Sectional Data, *Canadian Journal of Statistics*, 22, pp. 511-528.
- Roberts, C. (2007). Mixing modes of data collection in surveys: A methodological review. Review paper, NCRM/008, National Centre for Research Methods, City University London.
- Särndal, C.-E., and S. Lundström (2005). Estimation in Surveys with Nonresponse. New York: Wiley.
- Särndal, C.E., and B. Swensson (1987). A general view of estimation for two phases of selection with application to two-phase sampling and nonresponse. *International Statistical Review*, 55, pp. 279-294.
- Särndal, C-E., B. Swensson and J. Wretman (1992). *Model Assisted Survey Sampling*. New York: Springer Verlag.
- Scott, A.J. and T.M.F. Smith (1974). Analysis of Repeated Surveys using Time Series Methods. *Journal of the American Statistical Association*, 69, pp. 674-678.
- Scott, A.J. , T.M.F. Smith, and R.G. Jones (1977). The Application of Time Series Methods to the Analysis of Repeated Surveys. *International Statistical Review*, 45, pp. 13-28.
- Singh, A.C. (1996). Combining information in survey sampling by modified regression. *Proceedings of the Section on Survey Research Methods*, American Statistical Association, 120-129.
- Singh, A.C., B. Kennedy and S. Wu (2001). Regression Composite Estimation for the Canadian Labour Force Survey with a Rotating Panel Design. *Survey Methodology*, 27, pp. 33-44.

- Tam, S.M. (1987). Analysis of Repeated Surveys using a Dynamic Linear Model. *International Statistical Review*, 55, pp. 63-73.
- Tiller, R.B. (1992). Time Series Modelling of Sample Survey Data from the U.S. Current Population Survey. *Journal of Official Statistics*, 8, pp. 149-166.