



www.makswell.eu

Horizon 2020 - Research and Innovation Framework Programme

Call: H2020-SC6-CO-CREATION-2017

Coordination and support actions (Coordinating actions)

Grant Agreement Number 770643

Work Package 4

Time series and multivariate methodology for nowcasting well-being indicators and SDG's

Deliverable 4.1

Report on nowcasting and mixed frequency methods for the integrated analysis of well-being and SDG's

October 2019

Statistics Netherlands, Istat, Southampton University



This project has received funding from the European Union's Horizon 2020 research and innovation programme.



Deliverable D4.1

Report on the use of time series models for SDGs and well-being indicators

Authors

Statistics Netherlands: J.A. van den Brakel and C. Schiavoni

Istat: F. Bacchini, R. Iannaccone and D. Zurlo

Centre Camilo Dagum: I. Benedetti and T. Laureti

Southampton University: N. Tzavidis



Summary

The MAKSWELL project was set up to help strengthening the use of evidence and information on well-being and sustainability for policy-making in the EU, as also the political attention to well-being and sustainability indicators has been increasing in recent years.

To be written later.



1. Introduction	1
2. Data	3
2.1.Dutch LFS	3
2.2.Italian Well-being.....	6
2.2.1. Well-being indicators in the budget law	7
2.2.2. The Well-being indicators panel.....	8
2.3.Twitter.....	9
3. Methodology	11
3.1.Model for nowcasting unemployment.....	11
3.1.1. Time series model for monthly unemployment figures.....	11
3.1.2. Time series model for LFS data and claimant counts.....	12
3.1.3. Dynamic factor model for Google Trends and LFS time series	13
3.1.4. Estimation of structural time series models.....	15
3.1.5. Extensions of the dynamic factor model.....	15
3.2.Models for well-being	16
3.2.1. Dynamic factor model.....	16
3.2.2. Spatial panel model	16
4. Results	18
4.1.Dutch LFS	18
4.2.Panel	23
4.2.1. Dynamic factor model.....	23
4.2.2. Spatial panel model	24
4.3.Nowcasting of s80s20 index	24
4.4.Daily adjustment for the social mood	26
5. Discussion.....	29



1. Introduction

According to the Maxwell deliverable 1.1 (Tinto et al., 2018) we illustrated details for all the framework on well-being and SDG available across the European countries. We argued for an high degree of heterogeneity, both amid the indicators available and across countries and times. At the same time deliverable 2.1 presented several examples of non-traditional data sources.

This report aims to extend these results proposing a class of multivariate methods able to explore the well-being and SDG domain. The novelty of the results presented here is twofold: it has first to refer to the methodology, introducing dynamic factor models and spatial panel model, and secondly to look to application to well-being and SDG indicators such as unemployment, Italian well-being, subjective well-being, income inequality index and twitter.

Concerning multivariate methods, in Maxwell Deliverable 2.2 (van den Brakel et al., 2019) the area level model proposed by Fay and Herriot (1979) was proposed as potential small area prediction model using big data as covariates, since this model avoids the complex process of matching fussy big data sources with sample surveys at the unit level. It was also emphasized that most surveys conducted by national statistical institutes are conducted repeatedly over time. Therefore a natural approach for small area prediction as well as nowcasting is to extend the Fay-Herriot model with related information from previous editions of the survey. Accounts of regional small area estimation, where strength is borrowed over both time and space, include Rao and Yu (1994), Datta et al. (1999), You et al. (2003), You (2008), Pfeiffermann and Burck (1990), Pfeiffermann and Tiller (2006), Krieg and van den Brakel (2012), van den Brakel and Krieg (2016).

In this report multivariate structural time series models are developed to combine series obtained with repeated samples with related auxiliary series. This serves two purposes. First, extending the time series model an auxiliary series allows modelling the correlation between the unobserved components of the structural time series models, e.g. trend and seasonal components. If the model detects a strong correlation, then the accuracy of domain predictions will be further increased as illustrated by Harvey and Chung (2000) for the Labour Force Survey in the UK using a series of claimant counts.

Second, information derived from non-traditional data sources like Google trends or social media platforms are generally available at a higher frequency than series obtained with repeated surveys. This allows to use this time series modelling approach to make predictions for the survey outcomes in real time at the moment that the outcomes for the social media are available, but the survey data not yet. In this case the auxiliary series are used as a form of nowcasting (van den Brakel et al., 2017).

With a structural time series model a series is decomposed in a trend component, seasonal component, other cyclic components, regression component and an irregular component. For each component a stochastic model is assumed. This allows the trend, seasonal, and cyclic component but also the regression coefficients to be time dependent. If necessary ARMA components can be added to capture



the autocorrelation in the series beyond these structural components. See Harvey (1989) or Durbin and Koopman (2012) for details about structural time series modelling.

The time series observed with the repeated survey and the auxiliary series can both be expressed at the frequency of the survey. In that case the timeliness of the auxiliary series obtained with big data sources is utilized by making a first nowcast for the target parameter of the survey at the moment that the last observation of the auxiliary series becomes available and the survey estimate is still missing. Another way to exploit the timeliness of the auxiliary series is to define a mixed frequency model. In this case the model is expressed at the higher frequency, of the auxiliary series. This requires a disaggregation of the unobserved time series components of the target series to this higher frequency. After fitting the model, estimates for the survey parameters are obtained by aggregating the underlying components to a monthly frequency. Details of mixed frequency state-space models are described in Harvey (1989), Ch. 6.3, Durbin and Quenneville (1997), Moauro and Savio (2005).

With modern big data sources like Google Trends and Internet, a large amount of potential auxiliary series can be derived easily. Combining them directly in a multivariate structural time series models results in large models with many parameters, which on its turn result in models with reduced prediction power. This so called high dimensionality problem can be solved with dynamic factor models, which allows formulating parsimonious models despite a large amount of auxiliary series (Boivin and Ng, 2005, Stock and Watson, 2002a,b, Marcellino et al., 2003, Giannone et al., 2008, Doz et al., 2011). In section 3.1.3 such a dynamic factor model will be developed to combine time series observed with the Dutch LFS with a large amount of Google Trend series. The time series model for the LFS requires a complex set up in order to account for the rotating panel structure of this survey. This hampers the formulation of a dynamic factor model on the weekly frequency. Therefore nowcasting is in this paper performed by a model where all series are defined at the monthly frequency.

Dynamic factor models are explored also using the Italian well-being framework with the aim of exploring to correlation of the estimated factor in an environment characterized by an high number of three dimension: space (regions), time, and a consistent number of indicators. These three dimensions are also explicitly address by means of a space panel model.

With regards to nowcasting we present also the bridge model applied at Istat to nowcast the income inequality index using the preliminary estimation from the national accounts, that are timeless compared to the traditional picture related to the EU-silc survey.

Finally we present the main characteristics of the new sentiment indicator disseminated by Istat as an experimental statistics based on tweets.



2. Data

2.1. Dutch LFS

The objective of the Dutch LFS is to provide reliable information about the Dutch labour force. 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. Each month a stratified two-stage cluster design of addresses is selected. 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 and can be regarded as the ultimate sampling units. Most target parameters of the LFS concern people aged 15 through 64 years. Different subpopulations are oversampled to improve the accuracy of the official releases, for example addresses with persons registered at the employment office and subpopulations with low response rates.

Since October 1999, the LFS is conducted as a rotating panel design. Each month a new sample, drawn according the aforementioned sample design, enters the panel. Each sample is observed five times at quarterly intervals. The sample that is observed for the j -th time, is called the j -th wave of the panel, $j = 1, \dots, 5$. The sample that has been observed for the fifth time, leaves the panel. According to this rotation scheme, each month data are collected in five different waves. Data in the first wave are collected by means of computer assisted personal interviewing (CAPI). The respondents aged 15 through 64 years are re-interviewed four times 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. Participation of households with the Dutch LFS is on a voluntary basis. The monthly gross sample size for the first wave averaged about 8000 addresses commencing the moment that the LFS changed to a rotating panel design and gradually declined to about 6500 addresses in 2010. In July 2010 the data collection in the first wave changed from CAPI to a mix of CAPI and CATI. In 2012 the data collection in the first wave changed again to a sequential mixed mode design, starting with web interviewing and a follow up of CAPI and CATI. After these redesigns, the sample size increased to about 8000 addresses per month again. The response rate is about 55% in the first wave and in the subsequent waves about 90% with respect to the responding households from the previous wave.

Key parameters of the LFS are the employed, unemployed and total labour force, which are defined as population totals. Another important parameter is the unemployment rate, which is defined as the ratio of the unemployed labour force over the total labour force. Monthly estimates for these parameters are produced at the national level as well as a breakdown in six domains that is based on the cross classification of gender and three age classes.

There are two major problems with this survey. The first problem is that the monthly sample size of the LFS is too small to rely on the GREG estimator to produce timely official monthly statistics about



the employed and unemployed labour force. Therefore many national statistical institutes publish rolling quarterly figures about the labour force each month. Rolling quarterly figures have the obvious disadvantages that monthly seasonal patterns are smoothed out and that they are less timely since the monthly publications refer to the latest rolling quarter instead of the latest month.

The second problem is that there are substantial systematic differences between the subsequent panels due to mode and panel effects. This is a well-known problem for rotating panel designs, and in the literature this is referred to as rotation group bias (RGB), Bailer (1975). At the moment that the LFS changed from a cross-sectional survey to a rotating panel design in October 1999, the effects of the RGB on the outcomes of the LFS became very visible. This was the direct cause for developing procedures that account for this RGB.

Figure 2.1 illustrates the RGB for the unemployed labour force. The series of the GREG estimates of the first panel are compared with the average of the GREG estimates of the four subsequent panels. The GREG estimates for the unemployed labour force in the subsequent panels are systematically smaller than in the first panel. The RGB is a consequence of different non-sampling errors like selective non-response, panel attrition, mode-effects, effects due to differences between the CAPI questionnaire and the CATI questionnaire, and panel effects.

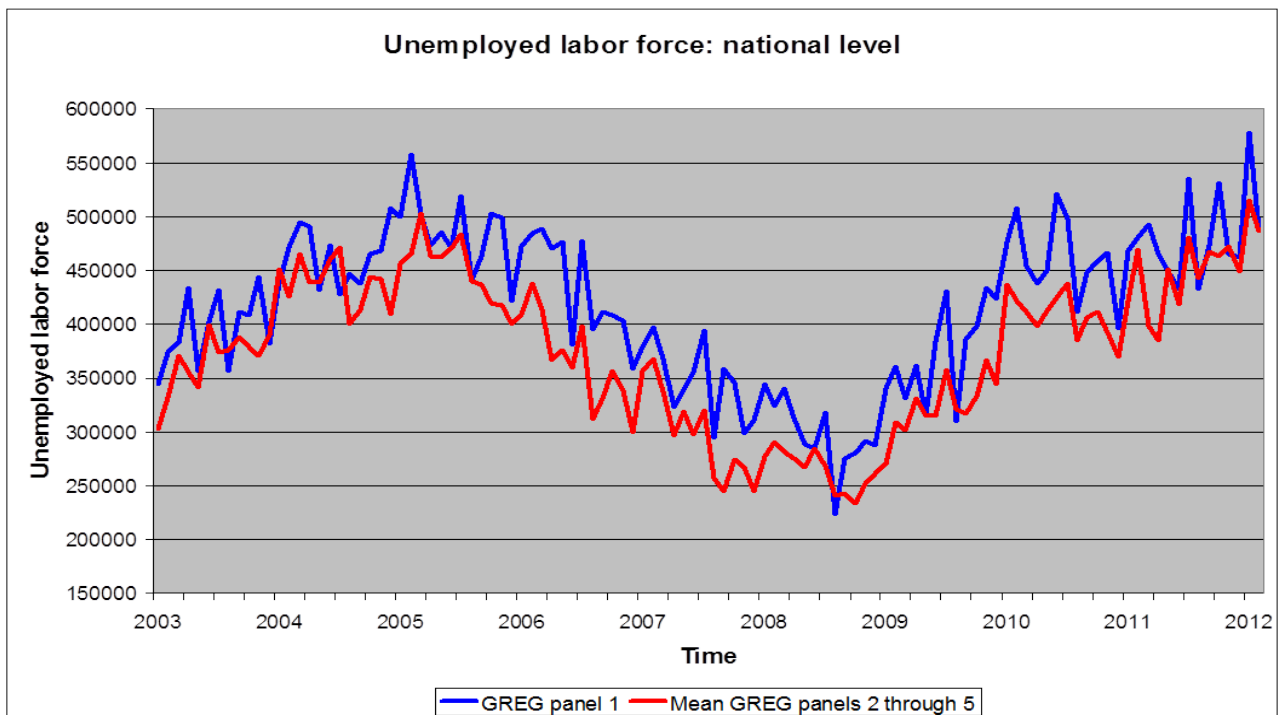


Figure 2.1: Comparison of estimates.

Problems with small sample sizes and RGB are solved with a five dimensional state-space model, initially proposed by Pfeffermann (1991). Monthly labour force figures are obtained with the following estimation procedure. As explained above, each month data are collected in five independent waves. The general regression GREG estimator is applied to produce five independent estimates for a target



parameter. Inclusion probabilities reflect the sampling design described above as well as the different response rates between geographic regions. The weighting scheme is based on a combination of different socio-demographic categorical variables. This results in five series of monthly GREG estimates for each target parameter, which are the input for a multivariate structural time series model described in Section 3.1.1. With this model reliable estimates for the population parameters are obtained by taking advantage of the sample information observed in previous periods. The model also accounts for RGB and autocorrelation induced by the rotating panel design. Since 2010 this approach is applied to produce official monthly figures about the labour force, (van den Brakel and Krieg, 2015).

The question addressed in this report is how the precision and timeliness of the monthly labour force figures can be further improved by taking advantage of additional auxiliary series. Two data sources are considered; claimant counts and Google trends.

The univariate auxiliary series of claimant counts represents the number of people claiming unemployment benefits. It is an administrative source, which is not available for every country, and, as for the Netherlands, it has the same publication delay of the labour force, i.e. both LFS estimates and the claimant counts for period t become available in $t + 1$. It is anticipated that this series will at least improve the precision of the time series model estimates for the monthly unemployment if we extend the state space model used by Statistics Netherlands in order to combine the survey data with this auxiliary series.

Besides claimant counts, the majority of the information related to unemployment is nowadays available on the internet; from job advertisements to resumé's templates and websites of recruitment agencies. We therefore follow the idea originating in Choi and Varian (2009) and Askitas and Zimmermann (2009) of using job-related terms searched on Google in the Netherlands. Since 2004, these time series are freely downloadable in real-time from the Google Trends tool, on a monthly or higher frequency. As from the onset it is unclear which search terms are relevant, and if so, to which extent, care must be taken not to model spurious relationships with regards to the labour force series of interest, which could have a detrimental effect on the estimation of unemployment, such as happened for the widely publicized case of Google Flu Trends (Lazer et al., 2014).

The Dutch labour force is subject to a one-month publication delay. In order to have more timely and precise estimates of the unemployment, we extend the model by including, respectively, auxiliary series about job search behaviour from weekly/monthly Google Trends and monthly claimant counts in the Netherlands.

Google Trends are indexes of search activity. Each index measures the fraction of queries that include the term in question in the chosen geography at a particular time, relative to the total number of queries at that time. The maximum value of the index is set to be 100. According to the length of the selected period, the data can be downloaded at either monthly, weekly, or higher frequencies. The series are standardized according to the chosen period and their values can therefore vary according to the period's length (Stephens-Davidowitz and Varian, 2015). We use weekly Google Trends for each search term, and are denoted x_t^{GT} .



Figure 2.2 displays the time series of the five waves of the unemployed labour force, together with the claimant counts and an example of job-related Google query. They all seem to be following the same trend, which already shows the potentiality of using this auxiliary information in estimating the unemployment.

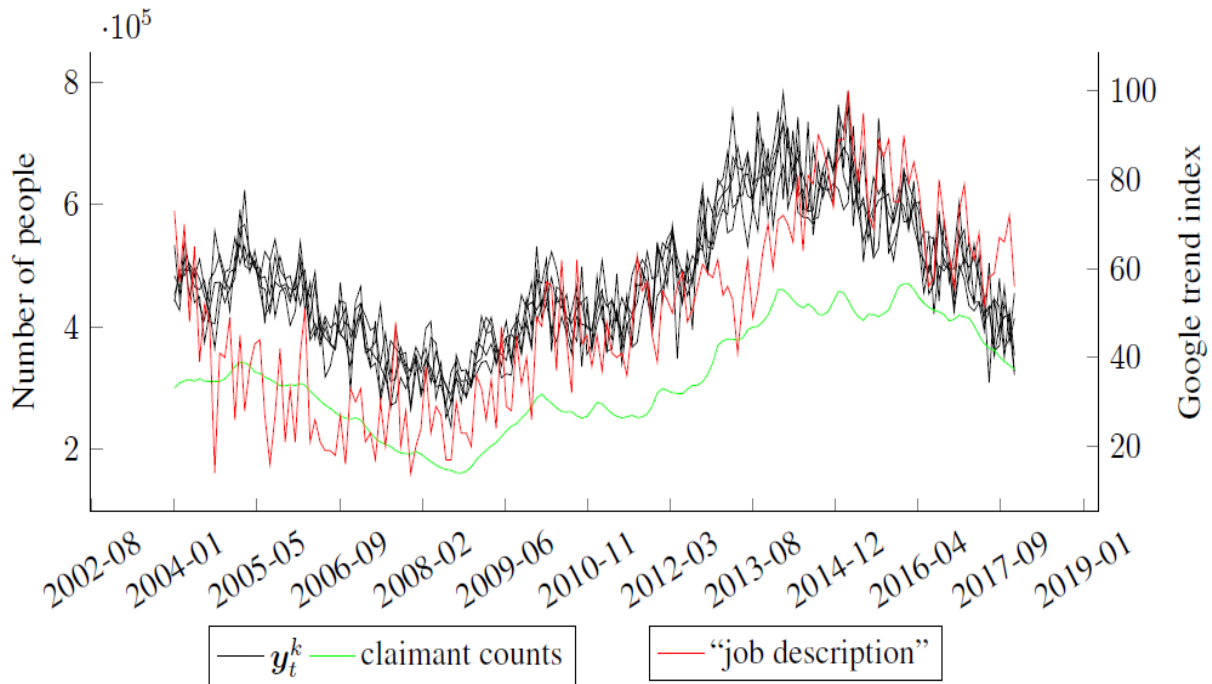


Figure 2.2: Comparison of LFS series, claimant counts and Google trend series for the search term ‘job description’.

2.2. Italian Well-being

The Italian National Institute of Statistics (Istat), together with the National Council for Economics and Labor (CNEL), launched in December 2010 an inter-institutional initiative aimed at developing a multi-dimensional approach for the measurement of “equitable and sustainable well-being” (Benessere equo e sostenibile), in line with the recommendations issued by the OECD and the Stiglitz Commission (see Stiglitz et al. 2009).

The 12 selected domains are divided into 2 typologies, 9 of them are defined as outcome domains and are those related to dimensions which have a direct impact on human and environmental well-being (?); the remaining 3 domains are defined as drivers of well-being, measuring functional elements to improve the well-being of the community and the surrounding environment ¹. The domains are:

- Outcome: health; education and training; work and life balance; economic well-being; social relationship; security; landscape and cultural heritage; environment; subjective well-being;
- Driver: politics and institutions; innovation, research and creativity; quality of services.

¹ This section relates to Bacchini et al. (2018)



In 2018 the importance attributed by citizens to each of the 12 domains of Bes in the individual perception of well-being was tested by a qualitative survey, which is an ideal update of that carried out in 2011 in the definition phase of the Bes domains.

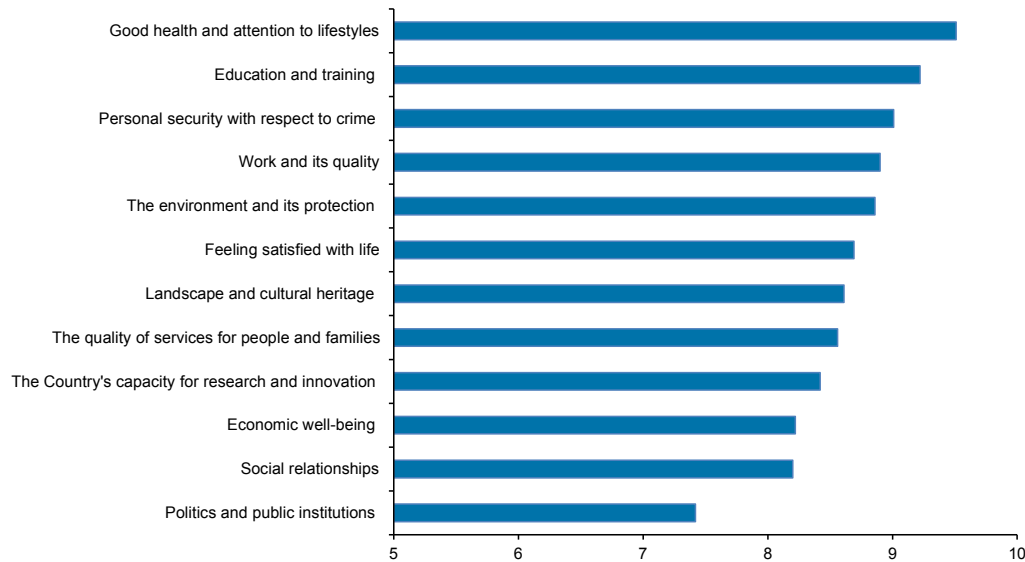


Figure 2.3: Average score attributed to the Bes domains (between 0 and 10). Italy. Year 2018.

In general, the 12 domains are confirmed relevant in defining the concept of well-being. Almost all of them receive an average rating of more than 8 (out of 10, Fig. 2.3). The only exception is the domain of Politics and institutions which received an average rating of 7.4, testifying a lower consideration from part of the citizens towards the different expressions of the public thing.

Very high scores, at least 9, are attributed to health, education and training, and personal safety, three cornerstones of individual well-being. The other Bes domains receive scores between 8 and 9, first of all the domain on work and quality of work, then gradually the others to end with economic well-being and social relationships (both 8.2).

The variability of the scores is however quite limited, with a substantial homogeneity of the evaluations expressed by different population groups.

The 12 domains, were originally populated with 134 indicators . However, the framework is considered as an open lab, and the set of indicators is reviewed annually to consider emerging information needs and methodologies. The Bes initiative has also been an important input to stimulate the production of new data on well-being. New questions were included in pre-existing surveys to be able to answer these needs. For instance, questions on trust in institutions and questions on perception of landscape and environment were added in the annual multipurpose survey on *Aspects of daily life*. According to this revision process the last edition of the annual report on Bes was based on 130 indicators.

2.2.1. Well-being indicators in the budget law

As we have seen in deliverable 1.1 (? Italy is one of the European country that consider well-being indicators directly in the budget law. For this scope, a scientific committee defines a selection of 12 indicators out of the 130 included in the Bes framework, namely:



Table 2.1: Number of indicators updated with 3-months time lag, by method

No.	Method	Source
7	Currently available	Istat, Ministry of Justice, Cresme
3	Ad hoc estimates on provisional data	Istat, Ministry of Interiors
2	Models for flash estimates	Istat and Istat based on Ispra data

1. Mean adjusted income (per capita)
2. Income inequality (quintile ratio)
3. Incidence of Absolute poverty;
4. Life expectancy in good health at birth
5. Overweight and obesity
6. Early school leavers
7. Non-participation in employment
8. Employment rate of women aged 25-49 with preschool children vs women without children
9. Victims of predatory crime
10. Mean length of civil justice trials
11. CO₂ and other greenhouse gas emissions (tons x inhab.)
12. Illegal Building

The adoption of this framework implies the alignment of data production to the law's provision, that in turn requires both to fasten the production process and to provide flash estimates based on provisional data, or even to implement forecasting models for those indicators whose data would be too late.

Table 2.2 shows the solutions adopted for the twelve indicators and the responsible body: the effort was a collective one, even though a great part of activities were borne by Istat, that had also a coordination role.

In the next section we present the flash estimation related to the income inequality.

2.2.2. The Well-being indicators panel

As we have reported, the complete dataset of Italian well-being refers to 130 indicators. Each of this indicator is observed both at national as well as for the 20 regions (Italy). The data available spans



Table 2.2: Number of indicators selected for the time series analysis

Domain Health	6
Education and training	5
Work and Life balance	6
Economic Well-Being	6
Social relationship	5
Politics and Institutions	2
Safety	4
Subjective Well-Being	1
Landscape and Cultural Heritage	3
Environment	5
Research and Innovation	2
Quality of services	4

from 2004 to 2017.

However due to an annual process of revision not all the indicators are available for all the years of the period. For our analysis we built up a balanced panel of 47 indicators available for the whole period for and for all Italian regions.

2.3. Twitter

Use of big data has starting to enter in the production process of the national institute of statistics. We have provided several example of that in deliverable 2.1 and 2.2. Sometimes the use of big data is related to new indicators based entirely on this new source. To emphasize this new indicators both Eurostat as well the national institute of statistics have create a new section on the own website called **Experimental statistics** (Figure 2.4).

Inside this framework, since october 2018 Istat has started to release an experimental index with quarterly frequency, based on Twitter data: the Social Mood on Economy Index (Figure 2.5

This new statistical instrument has been implemented to enable high-frequency (i.e. daily) measures of the Italian sentiment on the state of the economy. These measures are derived from samples of public tweets in Italian, which are captured on real time (see Zardetto (2018) for a general description of the methodology).

Twitter's streaming API is used to collect samples of public tweets matching a filter made up of 60 relevant keywords (actual words or phrases). A subset of these keywords has been borrowed from the questionnaire items of the Italian consumer confidence survey. To compute daily index values, all the tweets collected (about 50.000 every day) are firstly cleaned and normalized. For each tweet a sentiment analysis procedure is applied calculating positive and negative sentiment scores. For this purpose, message words are matched against entries of an Italian sentiment lexicon, namely a vocabulary whose lemmas are associated to pre-computed sentiment scores. Atomic scores of matched words are then averaged to yield tweet-level scores. Subsequently, tweets are clustered according to



EXPERIMENTAL STATISTICS

In line with Eurostat and other National Statistical Institutes, Istat experiments with the use of new sources and the application of innovative methods in producing data. The results are made available for users' use and evaluation.

They are experimental statistics because they do not respect all the steps necessary to test new methodologies, to transform them into technological and organisational solutions, to verify if quality requirements and harmonisation rules are fulfilled.

But their potential is really high; they fill knowledge gaps in a timely way; they serve as a driving force for new analyses and indicators; they guarantee an important information support to policies.

To ease users in finding and utilising them, experimental statistics produced by Istat are organised in four different typologies:

1. **Non-standard classifications** produced on the basis of the official taxonomies defined at an international level and currently used by Istat, or proposed as experimental within analysis and research activities based on microdata processing
2. **New indicators** produced through the integration of a multiplicity of official and non-official sources; in this case, the focus is on phenomena under investigation rather than on statistical sources used to describe them
3. **Interpretation frameworks and analysis** of complex phenomena obtained through the integration of official sources
4. Results of **Experiments on Big Data**, characterised, by their very nature, by the use of non-official sources.

NON-STANDARD CLASSIFICATIONS

NEW INDICATORS

INTERPRETATION FRAMEWORKS

EXPERIMENTS ON BIG DATA



FEEDBACK

You are invited to leave your observations, comments and suggestions on our experimental statistics. Please email to: statistiche-sperimentali@istat.it

Figure 2.4: Experimental statistics website at Istat

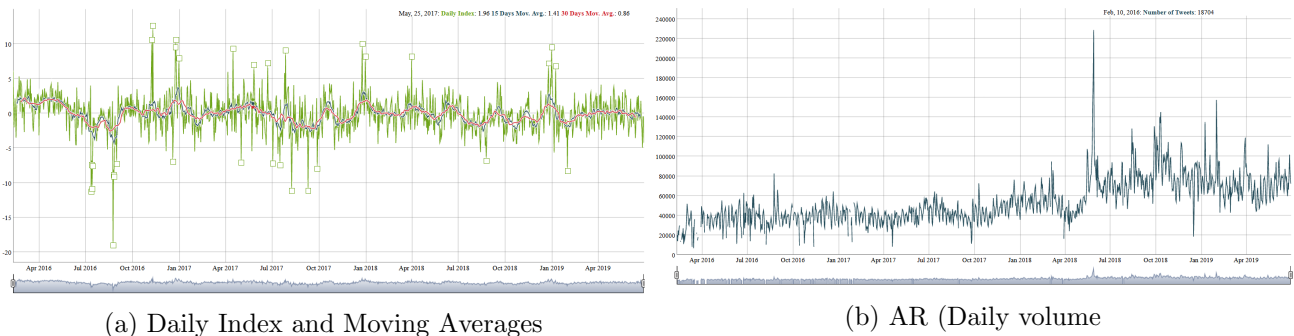


Figure 2.5: Social mood on economy index

their sentiment scores into three mutually exclusive classes: positive, negative and neutral tweets. The daily index value is derived as an appropriate central tendency measure of the score distribution of the tweets belonging to the positive and negative classes. As a last step, the daily index is linearly transformed in such a way that its long-run mean is zero. Special care has been devoted to make the index robust against possible contaminations by off-topic tweets that might pass the filter. To this end, a check system has been put in place, which periodically searches for anomalous values in the daily time series by means of two independent and complementary outlier detection routines.



3. Methodology

3.1. Model for nowcasting unemployment

3.1.1. Time series model for monthly unemployment figures

The data of the Dutch LFS and the way input series are calculated for the five dimensional time series model to produce monthly labour force figures is described in section 2.1. Let θ_t denote the unknown population total of the unemployed labour force in month t . Based on the rotating panel design, described in Section 2.1, data are collected in five different waves. Based on these data five independent GREG estimates for θ_t can be constructed. Let y_t^j denote the GREG estimate for the unknown population total θ_t based on the observation in wave j , $j = 1, \dots, 5$. This results in a five dimensional time series; $\mathbf{y}_t = (y_t^1, \dots, y_t^5)'$, for $t = 1, \dots, T$, which is the input for the following structural time series model (Pfeffermann, 1991, van den Brakel and Krieg, 2015)

$$\begin{pmatrix} y_t^1 \\ y_t^2 \\ \vdots \\ y_t^5 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \theta_t + \begin{pmatrix} 0 \\ \lambda_t^2 \\ \vdots \\ \lambda_t^5 \end{pmatrix} + \begin{pmatrix} e_t^1 \\ e_t^2 \\ \vdots \\ e_t^5 \end{pmatrix} \Leftrightarrow \mathbf{y}_t = \mathbf{1}_{[5]} \theta_t + \boldsymbol{\lambda}_t + \mathbf{e}_t. \quad (3.1)$$

The unknown population total is modelled in the first component of (3.1) with a smooth trend, say L_t , for the low frequency variation, a trigonometric seasonal component, say S_t , for the cyclic variation within one year and a white noise, say I_t , for the unexplained high-frequency variation, i.e.

$$\theta_t = L_t^\theta + S_t^\theta + I_t^\theta, \quad (3.2)$$

with the smooth trend model defined as:

$$\begin{aligned} L_t^\theta &= L_{t-1}^\theta + R_{t-1}^\theta, \\ R_t^\theta &= R_{t-1}^\theta + \eta_t^\theta \\ \eta_t^\theta &\simeq \mathcal{N}(0, \sigma_{\eta, \theta}^2) \end{aligned} \quad (3.3)$$

the trigonometric seasonal component defined as

$$\begin{aligned} S_t^\theta &= \sum_{j=1}^{J/2} \gamma_{j,t}^\theta \\ \gamma_{j,t}^\theta &= \gamma_{j,t-1}^\theta \cos\left(\frac{\pi j}{J/2}\right) + \gamma_{j,t-1}^{*,\theta} \sin\left(\frac{\pi j}{J/2}\right) + \omega_{j,t}^\theta \\ \gamma_{j,t}^{*,\theta} &= \gamma_{j,t-1}^{*,\theta} \cos\left(\frac{\pi j}{J/2}\right) - \gamma_{j,t-1}^\theta \sin\left(\frac{\pi j}{J/2}\right) + \omega_{j,t}^{*,\theta} \\ \omega_{j,t}^\theta &\simeq \mathcal{N}(0, \sigma_{\omega, \theta}^2) \quad \omega_{j,t}^{*,\theta} \simeq \mathcal{N}(0, \sigma_{\omega, \theta}^2) \\ j &= 1, \dots, J/2 \end{aligned} \quad (3.4)$$

and the white noise defined as $I_t^\theta \simeq \mathcal{N}(0, \sigma_{I, \theta}^2)$.



The second component of (3.1) models the RGB. Under the assumption that the observations obtained in the first wave are the most reliable and therefore unbiased, the relative bias between the follow-up waves and the first wave are modelled with random walks;

$$\begin{aligned}\lambda_t^j &= \lambda_{t-1}^j + \eta_{\lambda,t}^j, \quad j = 2, \dots, 5, \\ \eta_{\lambda,t}^j &\simeq \mathcal{N}(0, \sigma_{\eta,\lambda,j}^2).\end{aligned}\tag{3.5}$$

The third component of (3.1) models the autocorrelation in the survey errors in the follow-up waves, induced by the rotating panel design. In order to account for this autocorrelation, the survey errors are treated as state variables. To account for heteroscedasticity due to changing sample sizes over time, the state variables are scaled with the standard errors of the GREG estimates. This gives rise to the following model for the sampling errors

$$\begin{aligned}e_t^j &= c_t^j \tilde{e}_t^j, \quad c_t^j = \sqrt{\text{Var}(y_t^j)}, \quad j = 1, \dots, 5, \\ \tilde{e}_t^1 &\sim N(0, \sigma_{\nu_1}^2), \\ \tilde{e}_t^j &= \delta \tilde{e}_{t-3}^{j-1} + \nu_t^j, \quad \nu_t^j \sim N(0, \sigma_{\nu_j}^2), \quad j = 2, \dots, 5, \quad |\delta| < 1. \\ \text{Var}(\tilde{e}_{j,t}^k) &= \frac{\sigma_{\nu_j}^2}{(1 - \delta^2)}, \quad j = 2, \dots, 5.\end{aligned}\tag{3.6}$$

3.1.2. Time series model for LFS data and claimant counts

Let x_t^{CC} denote the claimant counts series which can be considered as an auxiliary series for the LFS time series model (3.1). One approach to improve the predictions for θ_t is to extend model (3.1) with a regression component; $\theta_t = L_t^\theta + S_t^\theta + \beta x_t^{CC} + I_t^\theta$. The major drawback of this approach is that the auxiliary series will partially explain the trend and seasonal effect in θ_t , leaving only a residual trend and seasonal effect for L_t^θ and S_t^θ . This hampers the estimation of a trend for the target variable. This problem is circumvented by modelling both series in a multivariate structural time series model:

$$\begin{pmatrix} \mathbf{y}_t \\ x_t^{CC} \end{pmatrix} = \begin{pmatrix} \mathbf{1}_{[5]}(L_t^\theta + S_t^\theta + I_t^\theta) \\ L_t^{CC} + S_t^{CC} + I_t^{CC} \end{pmatrix} + \begin{pmatrix} \boldsymbol{\lambda}_t \\ 0 \end{pmatrix} + \begin{pmatrix} \mathbf{e}_t \\ 0 \end{pmatrix}.\tag{3.7}$$

In (3.7) the LFS series and the claimant count series have their own trend, seasonal and disturbance term. In this application a smooth trend model is assumed for both series. The relation between both series can now be modelled via the covariance structure of the slope disturbances of both trend components, i.e.

$$\begin{aligned}L_t^z &= L_{t-1}^z + R_{t-1}^z, \\ R_t^z &= R_{t-1}^z + \eta_t^z, \quad z \in (\theta, CC), \\ \begin{pmatrix} \eta_t^\theta \\ \eta_t^{CC} \end{pmatrix} &\simeq \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{\eta,\theta}^2 & \rho_{\theta,CC}\sigma_{\eta,\theta}\sigma_{\eta,CC} \\ \rho_{\theta,CC}\sigma_{\eta,\theta}\sigma_{\eta,CC} & \sigma_{\eta,CC}^2 \end{pmatrix}\right).\end{aligned}$$



If the model detects a strong correlation between the trends of both series, then the trends of both series will develop into the same direction more or less simultaneously. In this case the additional information from the auxiliary series will result in an increased precision of the estimates of the target series θ_t . In the case of strong correlation between the disturbances of the trends, i.e. if $\rho_{\theta,CC} \rightarrow 1$, the trends are said to be cointegrated. This implies that the slope disturbances of both series simultaneously move up or down and that the slope disturbances of the auxiliary series can be perfectly predicted from slope disturbances of the target series. In that case there is one underlying common trend that drives the evolution of the trends of the two observed series. Cointegration increases the precision of the estimated trend and signal of the target series and allows for formulating more parsimonious models, which increases estimation efficiency. For a more detailed discussion about cointegration in the context of state space modelling, see Koopman et al. (2009), sections 6.4 and 9.1. The correlation between seasonal disturbance terms of both series can be modelled in a similar way.

3.1.3. Dynamic factor model for Google Trends and LFS time series

The Dutch labour force is subject to a one-month publication delay. In order to have more timely and precise estimates of the unemployment, we extend the model by including, respectively, auxiliary series about job search behaviour from weekly/monthly Google Trends.

Over the last decade, the number of non-traditional data sources that can be considered in the production of official statistics, is rapidly increasing. Particularly information derived from social media messages from Twitter and internet search behaviour from Google Trends easily result in a large number potential auxiliary series. Combining them in a full multivariate structural time series model as outlined in the previous subsection limits the degrees of freedom for model fitting. Due to the so-called "curse of dimensionality" prediction power of such models will be low. From this perspective, factor models are developed to formulate parsimonious models, despite a large number of auxiliary series are considered. Factor models are developed and widely applied by central banks to nowcast GDP on quarterly frequency using a large amount of related series observed on a monthly frequency (Boivin and Ng, 2005, Stock and Watson, 2002a,b, Marcellino et al., 2003). More recently, Giannone et al. (2008), Doz et al. (2011) proposed a state-space dynamic factor model. They propose a two-step estimator. In a first step a small amount of common factors are extracted from a large set of series using principal component analysis. In a second step, the common factors are combined with the target series in a state space model and are fitted using the Kalman filter.

Let \mathbf{x}_t^{GT} denote the vector with n auxiliary series derived from Google Trends, where n is large. In a first step a dynamic factor model is assumed for the auxiliary series

$$\mathbf{x}_t^{GT} = \mathbf{\Lambda} \mathbf{f}_t + \boldsymbol{\epsilon}_t, \quad (3.8)$$

$$\mathbf{f}_t = \mathbf{f}_{t-1} + \boldsymbol{\omega}_t, \quad (3.9)$$

with \mathbf{f}_t a p dimensional vector containing a small set of common factors that capture the major part of co-movements from the set of auxiliary series \mathbf{x}_t , where $p \ll n$. Furthermore $\mathbf{\Lambda}$ denotes a $n \times p$ dimensional matrix with factor loadings and $\boldsymbol{\epsilon}_t$ an n -vector containing variable specific shocks (idiosyncratic components). More over, the variance of the idiosyncratic shocks are defined as $\boldsymbol{\Psi} = \text{Var}(\boldsymbol{\epsilon}_t)$. If the series in \mathbf{x}_t are stationary, then \mathbf{f}_t can be estimated with the principal components on



\mathbf{x}_t and $\mathbf{\Lambda}$ through OLS. Generally potential auxiliary series will be non-stationary. If it is assumed that the series in \mathbf{x}_t are first order integrated (I(1)), then \mathbf{f}_t and $\mathbf{\Lambda}$ can be estimated with principal components analysis (PCA) on the differenced data $\mathbf{x}_t - \mathbf{x}_{t-1}$ (Bai, 2004). For identifiability reasons it is assumed that $Var(\boldsymbol{\omega}_t) = \mathbf{I}_p$.

The estimation steps proceed as follows. As outlined before, in a first step the factor loadings $\mathbf{\Lambda}$, the factors \mathbf{f}_t and the covariance matrix $\boldsymbol{\Psi}$ are estimated with PCA on the Google Trends observed on a weekly frequency. Then the Google Trends are averaged over the weeks within each month to obtain a series on a monthly frequency.

In a second step, the time series for the LFS, Google Trends and claimant counts are combined in a multivariate structural time series model:

$$\begin{pmatrix} \mathbf{y}_t \\ x_t^{CC} \\ \mathbf{x}_t^{GT} \end{pmatrix} = \begin{pmatrix} \mathbf{1}_{[5]}(L_t^\theta + S_t^\theta + I_t^\theta) \\ L_t^{CC} + S_t^{CC} + I_t^{CC} \\ \hat{\mathbf{\Lambda}}\mathbf{f}_t \end{pmatrix} + \begin{pmatrix} \boldsymbol{\lambda}_t \\ 0 \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{e}_t \\ 0 \\ \boldsymbol{\epsilon}_t \end{pmatrix}. \quad (3.10)$$

In model (3.10), the estimated factor loadings $\hat{\mathbf{\Lambda}}$ and the estimated covariance matrix $\hat{\boldsymbol{\Psi}}$, obtained in the first step, are kept fixed. The factor loadings are re-estimated with the Kalman filter, assuming the following relations between the LFS trend, claimant counts trend and the Google Trend common factors:

$$\begin{aligned} L_t^z &= L_{t-1}^z + R_{t-1}^z, \\ R_t^z &= R_{t-1}^z + \eta_t^z, z \in (\theta, CC), \\ \mathbf{f}_t &= \mathbf{f}_{t-1} + \boldsymbol{\omega}_t, \\ \begin{pmatrix} \eta_t^\theta \\ \eta_t^{CC} \\ \boldsymbol{\omega}_t \end{pmatrix} &\simeq \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \\ \mathbf{0}_p \end{pmatrix}, \begin{pmatrix} \sigma_{\eta,\theta}^2 & \rho_{\theta,CC}\sigma_{\eta,\theta}\sigma_{\eta,CC} & \rho_{\theta,G1}\sigma_{\eta,\theta} & \dots & \rho_{\theta,Gp}\sigma_{\eta,\theta} \\ \rho_{\theta,CC}\sigma_{\eta,\theta}\sigma_{\eta,CC} & \sigma_{\eta,CC}^2 & 0 & \dots & 0 \\ \rho_{\theta,G1}\sigma_{\eta,\theta} & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_{\theta,Gp}\sigma_{\eta,\theta} & 0 & 0 & \dots & 1 \end{pmatrix} \right). \end{aligned}$$

The estimate of $\mathbf{\Lambda}$ is used in the design matrix of the measurement equation of the state space model since its knowledge is needed in order to apply the Kalman filter. We fix $\hat{\boldsymbol{\Psi}}$ because it is a high-dimensional covariance matrix and we could incur in the curse of dimensionality if we would have to estimate it by maximum likelihood. Restricting the sample covariance matrix of the idiosyncratic components as being diagonal is standard in the literature. The consistency of this two-step estimator has been originally proven in the stationary framework by Doz et al. (2011), and extended to the nonstationary case by Barigozzi and Luciani (2017).

The Kalman filter (second step) is applied to the whole state space model to re-estimate \mathbf{f}_t and to nowcast the variable of interest, L_t^θ , and $L_t^\theta + S_t^\theta$ providing unemployment estimates in real time before LFS data become available. Since in each week we can aggregate the weekly Google Trends to the monthly frequency, we can use the information available throughout the month to update the estimated factors and loadings of the auxiliary series. If the correlations between the factors and the



trend's slope of the target variable are large, this update should provide a more precise nowcast of L_t^θ , and $L_t^\theta + S_t^\theta$.

3.1.4. Estimation of structural time series models

A widely applied approach to fit structural time series models is to write them in state-space form and analyses them with the Kalman filter. The Kalman filter is a recursive procedure that runs from period $t = 1$ to T and gives, for each time period, an optimal estimate for the state variables based on the information available up to and including period t . These estimates are referred to as the filtered estimates. The filtered estimates of past state vectors can be updated, if new data after period t become available. This procedure is referred to as smoothing and results in smoothed estimates that are based on the complete time series. The Kalman filter assumes that the hyperparameters, i.e. the variance components of the stochastic processes for the state variables are known. This is generally not the case. In practice maximum likelihood estimates for the hyperparameters are obtained using a numerical optimization procedure (BFGS algorithm). Expressions for the state space representation of structural time series models and details of the Kalman filter can be found in Durbin and Koopman (2012). The state-space representation of the models from Subsections 3.1.1, 3.1.2, and 3.1.3 are given in Schiavoni et al. (2019).

Several software packages are available to fit structural time series models. Most standard structural time series models can be fitted with STAMP (Koopman et al., 2009). For more advanced models, more advanced software is required. One option is to implement these models in OxMetrics in combination with the subroutines of SsfPack 3.0, see Doornik (2009) and Koopman et al. (2009, 2008). The models fitted in this paper are implemented in R (Team, 2017).

3.1.5. Extensions of the dynamic factor model

In an attempt to extract more information from the Google Trend series, three extensions are explored.

The first approach is to target the Google Trends by using the Elastic Net as a variable selection procedure for the most important Google Trends to improve their forecast performance. We follow the approach proposed by Bai and Ng (2008) and regress the differenced estimated change in unemployment from the labour force model without auxiliary series, $\Delta \hat{R}_t^{k,y}$, on the differenced Google Trends using the penalized regression proposed by Hastie and Zou (2005). The tuning parameter for the penalty in the Elastic Net (λ) is chosen in order to minimize the Akaike Information Criterion (AIC) for a grid of values for the weight for the L1 and L2 norm in the range $[0.05, \dots, 0.95]$.

A second possible improvement is to include the lags of the Google Trends' factors. It is reasonable to assume that people might start looking for a job before becoming unemployed. We therefore propose a parsimonious method to let the innovation of the change in unemployment depend on the lags of the Google Trends' factor. Assume only one relevant factor for the Google Trends and consider a regression of $\eta_{R,t}^{k,y}$ on the lags of the differenced factor:

$$\eta_t^\theta = \sum_{j=1}^q \kappa_j u_{t-j} + w_t = \kappa_1 f_{t-1} + \sum_{j=2}^q (\kappa_j - \kappa_{j-1}) f_{t-j} - \kappa_q f_{t-q-1} + w_t^k, \quad w_t \sim N(0, \sigma_w^2).$$



η_t^θ is estimated from the labour force model without auxiliary series, and regressed on \hat{u}_t^k , estimated by PCA, in order to obtain ordinary least squares estimates of the parameters κ . The choice of the number of lags to be included in the regression is chosen by the AIC.

The third extension is an attempt to model the seasonality or cycle of the Google Trends' factors. The Google Trends' factors might capture cycles or the seasonality of the job search terms. Assume again only one relevant factor for the Google Trends. In line with Alonso et al. (2011), instead of deseasonalizing the Google Trends, we model the seasonality or cycle of the factors. We assume that f_t follows a seasonal ARIMA model. Once f_t is estimated by PCA, the parameters of the seasonal ARIMA model can be estimated by ordinary least squares and plugged in the transition equation of the state space model.

3.2. Models for well-being

3.2.1. Dynamic factor model

As pointed out by Forni et al. (2018): large-Dimensional Dynamic Factor Models represent each variable in the dataset as decomposed into a common component, driven by a small (as compared to the number of series in the dataset) and fixed (as the number of series grows) number of common factors and an idiosyncratic component. For the application to the Italian well-being indicators we consider the general representation provided in ?:

$$x_{it} = \lambda_{i1}F_{1t} + \lambda_{i2}F_{2t} + \dots + \lambda_{irt}F_{rt} + \xi_{it} \quad (3.11)$$

where the factors F_{it} and the loadings λ_{1t} represent the common component driven by a small and fixed number of common factors. The idiosyncratic component are assumed to be orthogonal across different variables or only weakly correlated, so that the covariance of the variables is mostly accounted for by the common components. In the proposed factor model, the factors and the loadings are estimated using the first r standard principal components.

Dynamic factor model have not so often applied to the well-being and SGD set of indicators due to their limits in terms of time span availability. However using the database of the Italian well-being framework we have been able to provide a first application that could be important both for the nowcasting of the indicators as well as to investigate in the differences amid the common components estimated across the domains and the regions.

3.2.2. Spatial panel model

Spatial panel data model capture spatial interaction across spatial units and over time. The general static panel model which include a spatial lag of the dependent variable and spatial autoregressive disturbance is the follow:

$$y = \lambda(I_T \otimes W_n)y + X\beta + u \quad (3.12)$$

where y is an $N_T \times 1$ vector of observations on the dependent variable. X is a $N_T \times k$ matrix of observations on the non-stochastic exogenous regressors. I_T is an identity matrix of dimension T , W_n is the $N \times N$ spatial weights matrix of known constants whose diagonal elements are set to zero, λ is the



spatial parameter.

The disturbance vector is the sum of two terms:

$$u = (i_T \otimes I_N)\mu + \epsilon \quad (3.13)$$

where u is a vector of time invariant individual specific effects (not spatially autocorrelated) and ϵ is a vector of spatially autocorrelated innovation that follow a spatial autoregressive process:

$$\epsilon = \rho(I_T \otimes W_n)e + \nu \quad (3.14)$$

with $\rho(|\rho| < 1)$ as the spatial autoregressive parameter, and $\nu_{it} \sim IID(0, \sigma_\nu^2)$ and $\epsilon_{it} \sim IID(0, \sigma_\epsilon^2)$.

In random effect model, the unobserved individual effects are uncorrelated with the other explanatory variables in the model. In this case $\mu_i \sim IID(0, \sigma_\epsilon^2)$.

The error term can be rewritten as:

$$\epsilon = \rho(I_T \otimes B_N^{-1})\nu \quad (3.15)$$

where $B_N = (I_N - \rho W_N)$. The error terms will be:

$$u = (i_T \otimes I_N)\mu + (I_T \otimes B_N^{-1})\nu \quad (3.16)$$

The use of the spatial panel model is implemented together with the testing for unit roots in time series (Levin et al. (2002)) where the null hypothesis is that each individual time series contains a unit-root against the alternative that each time series is stationary.



4. Results

4.1. Dutch LFS

As explained in Subsection 3.1.3, the Google series used in the model must be $I(1)$. It is therefore tested whether the Google Trends are non-stationary with the Elliott et al. (1996) augmented Dickey-Fuller (ADF) test, including a constant and a linear trend. To control for the overall significance level a moving block bootstrap approach that accounts for time and cross-sectional dependence among the units in the panel, proposed by Moon and Perron (2012), is applied. The Google Trends that resulted as being $I(1)$ from this multiple hypotheses test are used in the dynamic factor model. Whenever applying PCA or the Elastic Net, the Google Trends are first differenced and standardized.

Four different models are compared:

1. The labour force model (3.1) without auxiliary series as it is currently used in the production official monthly figures (baseline)
2. The labour force model with auxiliary series of claimant counts (CC) defined in (3.7),
3. The labour force model with auxiliary series of Google Trends (GT) defined in (3.10) but without the CC component
4. The labour force model with auxiliary series of claimant counts and Google Trends (CC & GT) defined in (3.10)

The latter three models are compared to the baseline model with an in-sample and an out-of-sample exercise. The period considered for the estimation starts in January 2004 and ends in December 2017 ($T = 167$ months). The out-of-sample nowcasts are conducted in real time (concurrently) in the last four years of the sample. Each week the model is re-estimated assuming that the current observations for the unemployed labour force and the claimant counts are missing.

The estimation uncertainty is measured as the average over the variance of the filtered state variables for the trend and slope of the LFS ($\hat{L}_{t|t}^\theta, \hat{R}_{t|t}^\theta$) and the signal ($\hat{\theta}_{t|t} = \hat{L}_{t|t}^\theta + \hat{S}_{t|t}^\theta$), i.e.

$$M\hat{S}E(\hat{\mathcal{Z}}) = \frac{1}{T-d} \sum_{t=d+1}^T \hat{V}ar(\mathcal{Z}_{t|t}), \quad \mathcal{Z}_{t|t} \in (\hat{L}_{t|t}^\theta, \hat{R}_{t|t}^\theta, \hat{\theta}_{t|t}) \quad (4.1)$$

where $\hat{V}ar(\mathcal{Z}_{t|t})$ is obtained from covariance matrix of the filtered state variables of the Kalman filter recursion.

The nowcast accuracy is measured with the variance of the one step ahead forecast for the trend and slope of the LFS ($\hat{L}_{t|t-1}^\theta, \hat{R}_{t|t-1}^\theta$) and the signal ($\hat{\theta}_{t|t-1} = \hat{L}_{t|t-1}^\theta + \hat{S}_{t|t-1}^\theta$), at the moment that the Google



Trends become available in respectively week $j = 0, 1, \dots, 5$ of month t , and the observations for the LFS and CC are still missing. For each month the Maximum Likelihood estimates are recalculated using the series observed until that time period. This means that the nowcasts analysis is done in real time. More formally, the nowcast uncertainty is defined as:

$$M\hat{S}E(\hat{\mathcal{Z}}) = \frac{1}{H} \sum_{t=T-H+1}^T \hat{V}ar(\mathcal{Z}_{t|t-1}), \quad \mathcal{Z}_{t|t-1} \in (\hat{L}_{t|t-1}^{\theta}, \hat{R}_{t|t-1}^{\theta}, \hat{\theta}_{t|t-1}) \quad (4.2)$$

Results for MSE's are reported as the relative $M\hat{S}E$ with respect to the baseline model. Values lower than one are therefore in favour of the model under consideration. Likelihood ratio (LR) tests are conducted to assess whether the correlation parameters are different from zero, and hence adding the auxiliary information might yield a significant improvement from the baseline model. Namely, the null hypotheses of the test for the CC, GT and CC & GT models are, respectively: $\rho_{CC} = 0$, $\rho_{GT} = 0$ and $\rho_{CC} = \rho_{GT} = 0$. The test statistics should be compared to the critical values of a χ^2 distribution with degrees of freedom equal to the number of parameters that are being tested.

Based on the PCA applied to the Google Trends, finally two common factors were selected for the dynamic factor model. Table 4.1 reports the maximum likelihood estimates for the hyper parameters for the four models. Table 4.2 reports the relative measures of in and out-of-sample performance when the monthly Google Trends are used.

Hyperparameter	discription	Baseline	CC	GT	CC & GT
$\hat{\sigma}_{\eta,\theta}$	slope LFS	2201.140	3023.983	2473.391	3096.768
$\hat{\sigma}_{\omega,\theta}$	seasonal LFS	0.020	0.020	0.020	0.020
$\hat{\sigma}_{\eta,\lambda}$	RGB LFS	1166.055	1214.057	1037.313	1114.421
$\hat{\sigma}_{v_1}$	sampling error wave 1	1.165	1.169	1.164	1.169
$\hat{\sigma}_{v_2}$	sampling error wave 2	1.139	1.139	1.143	1.141
$\hat{\sigma}_{v_3}$	sampling error wave 3	1.082	1.077	1.087	1.078
$\hat{\sigma}_{v_4}$	sampling error wave 4	1.128	1.144	1.135	1.144
$\hat{\sigma}_{v_5}$	sampling error wave 5	1.100	1.107	1.102	1.108
$\hat{\sigma}_{\eta,CC}$	slope CC		3606.933		3595.186
$\hat{\sigma}_{\omega,CC}$	seasonal CC		0.020		0.020
$\hat{\sigma}_{\epsilon,CC}$	white noise CC		1120.032		1.052
$\hat{\rho}_{\theta,CC}$	correlation LFS-CC		0.902		0.903
$\hat{\rho}_{\theta,G_1}$	correlation LFS-GT1			0.428	-0.038
$\hat{\rho}_{\theta,G_2}$	correlation LFS-GT2			-0.397	0.051
p-values from the LR test					
$H_0 : \hat{\rho}_{\theta,CC} = 0$			0.0004		0.0007
$H_0 : \hat{\rho}_{\theta,G_1} = 0$				0.3897	1.00
$H_0 : \hat{\rho}_{\theta,G_2} = 0$				0.3125	1.00

Table 4.1: Maximum likelihood estimate hyperparameters for the labour force model with auxiliary series of claimant counts and aggregated weekly Google Trends to the monthly frequency.

The estimated correlation with the claimant counts is large, more than 0.9, and remains such when including the Google Trends. The correlation with the two Google Trend factor levels is around 0.4



and -0.4 in the model with Google Trend factors as auxiliary series only. Note that the sign of the correlation with the factor is not relevant, as the factor, in PCA, is identified up to a sign. We are therefore only interested in the magnitude of this parameter. In the full model where claimant counts and the Google Trends are used both as auxiliary series, the strong correlation of the claimant counts remains while the correlation of the Google Trend factors shrinks to zero. The LR tests suggest that only the correlation between the slope disturbance terms of the LFS and the claimant counts is significantly different from zero, indicating a preference for the model which contains this auxiliary information rather than the Google Trends.

MSE	CC	GT	CC & GT
$\hat{MSE}(\hat{L}_{t t}^\theta)$	0.869	0.967	0.869
$\hat{MSE}(\hat{R}_{t t}^\theta)$	0.956	0.893	0.988
$\hat{MSE}(\hat{\theta}_{t t}^\theta)$	0.890	0.977	0.889
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$	0.715		
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$	0.929		
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$	0.729		
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$ after week 0		0.949	0.731
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$ after week 1		0.949	0.737
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$ after week 2		0.949	0.722
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$ after week 3		0.951	0.731
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$ after week 4		0.952	0.734
$\hat{MSE}(\hat{L}_{t t-1}^\theta)$ after week 5		0.938	0.703
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$ after week 0		0.882	0.930
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$ after week 1		0.881	0.924
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$ after week 2		0.881	0.941
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$ after week 3		0.883	0.936
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$ after week 4		0.887	0.922
$\hat{MSE}(\hat{R}_{t t-1}^\theta)$ after week 5		0.873	0.927
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$ after week 0		0.953	0.743
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$ after week 1		0.953	0.749
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$ after week 2		0.953	0.735
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$ after week 3		0.955	0.744
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$ after week 4		0.956	0.756
$\hat{MSE}(\hat{\theta}_{t t-1}^\theta)$ after week 5		0.943	0.717

Table 4.2: Estimation and nowcast MSE results for the labour force model with auxiliary series of claimant counts and aggregated weekly Google Trends to the monthly frequency relative to the baseline model.

The claimant counts improve the estimation accuracy of the trend and signal with 13% and 11% respectively. Although the claimant counts for period t become available in $t + 1$, they also improve the nowcast accuracy for the trend with almost 30% and the signal of 27%. The model with Google Trends gives a small improvement of the estimation accuracy and nowcast accuracy. The accuracy of the nowcast does not monotonically improve with the number of weeks. The full model both improve



the estimation accuracy as well as the nowcast accuracy, but this contribution comes mainly from the claimant counts, since the correlation between the slope disturbance terms of the LFS and the level disturbance terms of the Google Trend common factors shrinks to zero.

Concerning the proposed extensions of the dynamic factor model in Subsection 3.1.5 the following conclusions can be reported. Targeting the Google Trends with the Elastic Net does not increase the value of the correlation parameter, nor improves it the estimation and nowcast accuracy.

When including the lags of the factor, only one lag is always chosen by the AIC in each recursion of the out-of-sample exercise, and its estimated parameter is insignificant. The AIC would actually prefer to not include any lag, but we force at least one lag to be included in order to have a different model. This already suggests that the inclusion of the lag should not improve the accuracy of the estimation. For this extension, the correlation with the Google Trends' factors becomes really large, above 0.9, such that the LR test even rejects the null hypothesis of this parameter being equal to zero, but the measures of estimation and nowcast accuracy worsen with respect to the baseline model. The reason of these results might be found in the negative and significant first order autocorrelation of both \hat{u}_t . This means that η_t^θ is forced to be more correlated with \hat{u}_t because of its dependence on \hat{u}_{t-1} . Therefore, this model is miss-specified for u_t .

The estimated factor shows a cyclical pattern, especially when the factor loadings are estimated on the weekly Google Trends. Since the seasonal ARIMA model is fitted on the estimated monthly factor by PCA, $s = 12$. The number of lags and MA components are again chosen according to the AIC, which suggests an ARIMA(3, 1, 1) as the best model. This result indicates that the regular pattern in the factor does not repeat yearly, but rather quarterly. This model yields a slight improvement of the estimation and nowcast accuracy but the LR tests are not in favor of the model that includes Google Trend factors.

Figures 4.1, 4.2 and 4.3 display the nowcast of the slope of the trend, its level and the population parameter, of the four models. Especially from the first graph, it is evident that the models including the claimant counts slightly deviate from the baseline model. On the contrary, the baseline model and the model with Google Trends give similar results. The reason is likely due to the large correlation between the slope disturbance terms between the LFS and the claimant counts series. This is not the case for the GT model as the correlation with the search terms' factor is not significantly different from zero.

The assumptions of no serial correlation, heteroscedasticity and normality made throughout the paper can be tested on the standardized one-step ahead forecast errors. The test for autocorrelation is conducted with the Ljung-Box test for 4, 8, 12 and 16 lags. The Ljung-Box test only detects significant autocorrelation in the one-step ahead forecasts of the claimant counts. An F-test for heteroscedasticity does not detect heteroscedasticity in the standardized one-step ahead forecast errors of the LFS series and claimant count series. The Sharipo-Wilk test for univariate normality did not detect deviations from the standard normal distribution of the standardized one-step ahead forecast errors.

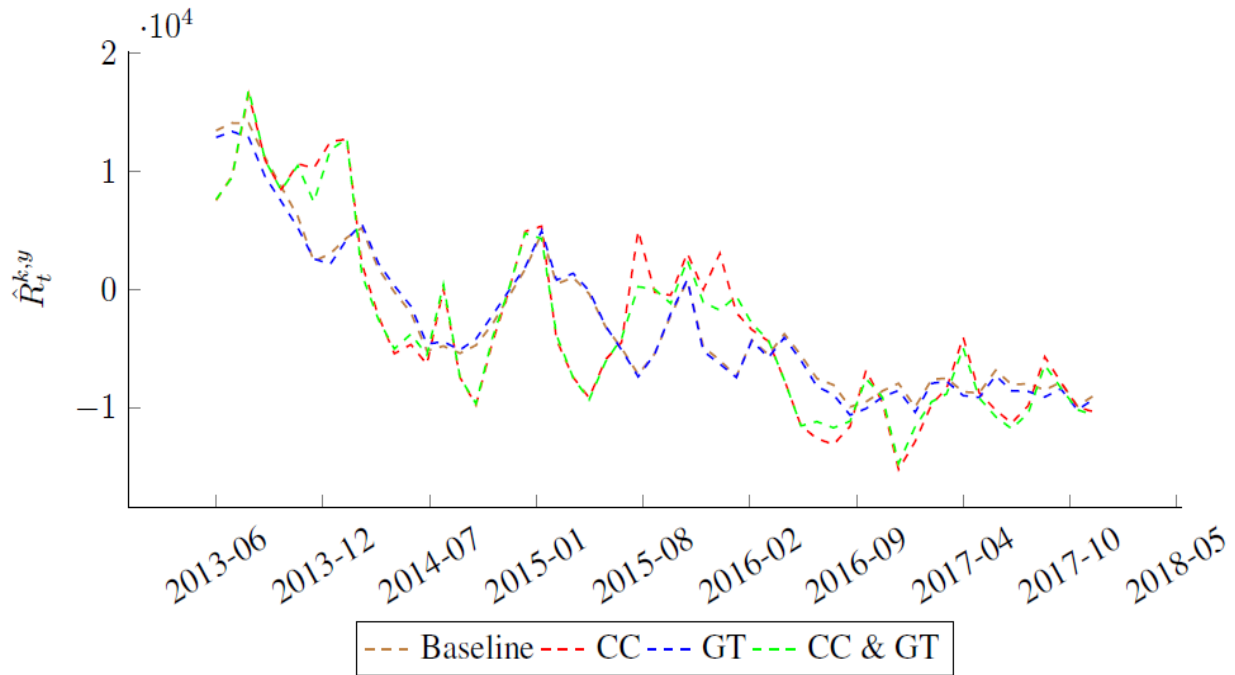


Figure 4.1: Nowcast of R_t for monthly unemployment under four different models.

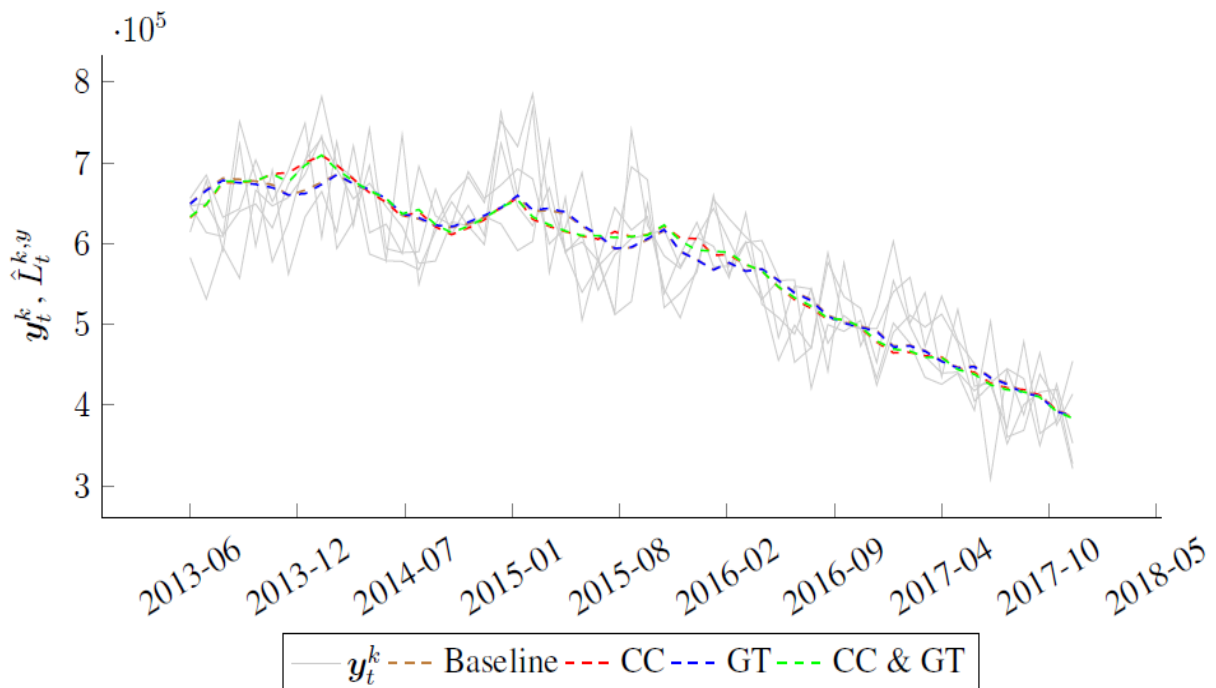


Figure 4.2: Nowcast of L_t for monthly unemployment under four different models.

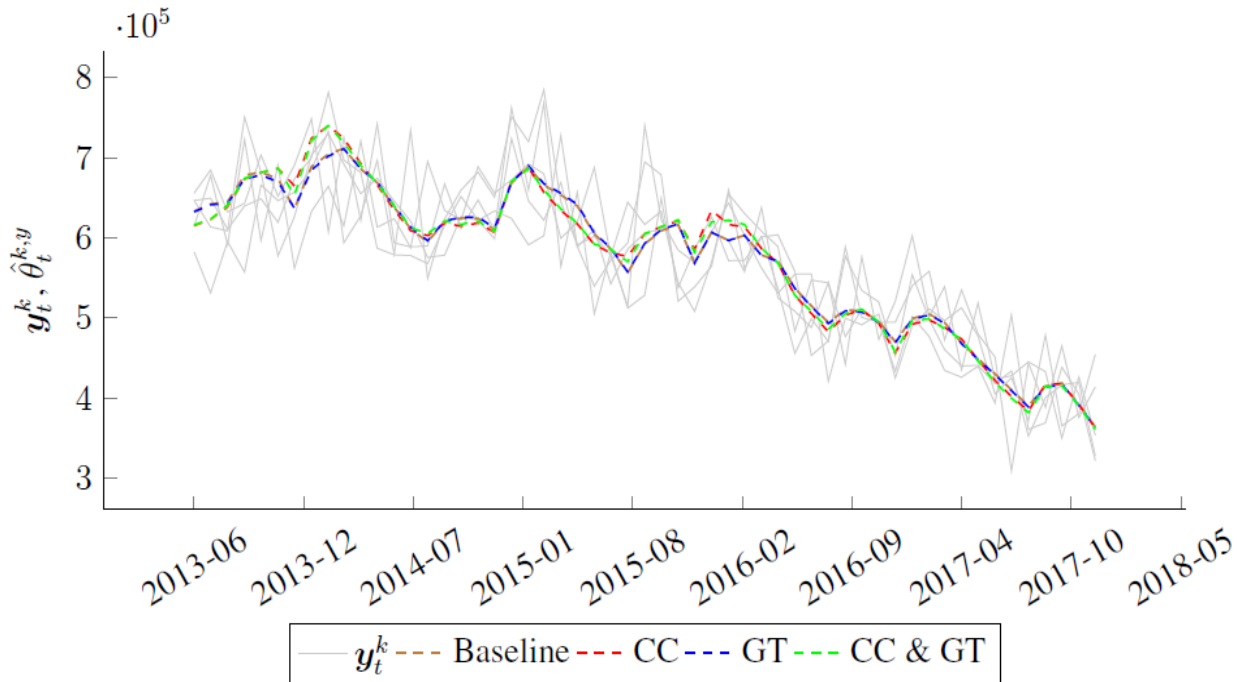


Figure 4.3: Nowcast of θ_t for monthly unemployment under four different models.

4.2. Panel

4.2.1. Dynamic factor model

We first refer to the application of the dynamic factor model for each regions using the set of 47 indicators span from 2004 to 2017. We first explore the correlation amide the factors estimated for each region. Figure 4.4 illustrates the correlation for the first 2 factors estimated for Italy, Piedmont and Sicily while figure 4.5 presents the evolution of the cators estimated along the time.

Factor 1	Italy	Piedmont	Sicily	Factor 2	Italy	Piedmont	Sicily
Italy	1,00	-0,2678	-0,1642	Italy	1,00	0,1688	-0,0006
Piedmont	-0,2678	1,00	-0,3908	Piedmont	-0,1688	1,00	0,1094
Sicily	-0,1642	0,3908	1,00	Sicily	-0,0006	0,1094	1,00

(a) Factor 1

(b) Factor 2

Figure 4.4: Correlation amid the factors and regions

Both figures could be associate with different interpretation of the common characteristics across the territory. The main drivers for Italy, Piedmont and Sicily are respectely:

- Satisfaction with family relations and friendship, Irregularities in electric power distribution,



Factor 1	Italy	Piedmont	Sicily	Factor 2	Italy	Piedmont	Sicily
Italy	1,00	-0,2678	-0,1642	Italy	1,00	0,1688	-0,0006
Piedmont	-0,2678	1,00	-0,3908	Piedmont	-0,1688	1,00	0,1094
Sicily	-0,1642	0,3908	1,00	Sicily	-0,0006	0,1094	1,00

(a) Factor 1

(b) Factor 2

Figure 4.5: Correlation amid the factors and regions

Separate collection of municipal waste

- Nutrition, Obesity, Satisfaction with means of transport
- Sedentarioussness, Impact of forest fires, Association funding

4.2.2. Spatial panel model

Starting from the dababase of the Italian well-being the application of the spatial panel model (SPM) refers to a subset of 57 indicators available for Italian regions for the periodo 2011-2017. For this application we use as dependent variable the Subjective Well-Being expressed as percentage of people aged 14 and over with a level of life satisfaction from 8 to 10 on total population aged 14 and over while independent variable for each domain are present in Figure 4.6

Subjective Well-Being	Percentage of people with high level of Life Satisfaction	dependent variable
Domain	Indicators	
Health	Life expectancy at birth (Bérenger and Verdier-Chouchane, 2007)	
Education and Training	Early leavers from education and training, (Skorikov, V. 2007)	
Work and Life balance	Unemployment rate (Alesina et al., 2014)	
Economic Well-Being	Economic inequalities (Clark, Fritters, and Shields 2007)	
Social relationship	Personal networks, voluntary activity and dimensions of social capital (Scrivens and Smith, 2014)	
Politics and Institutions	Trust in public institutions (police, legal system and government (Helliwell and Putnam, 2004; Hudson, 2006)	
Landscape and cultural heritage	Illegal building rate (D'Amato and Zoli, 2011)	
Research and Innovation	Impact of knowledge workers on employment (Engelbrecht, H. J., 2012)	
Safety	Fear of crime (Ferrer-i-Carbonell & Gowdy, 2007)	
Environment	Energy from renewable resources (Welsch, H., & Biermann, P., 2014).	

Figure 4.6: X variables used in the Space panel model

Estimation of the space model are reported in Figure 4.7. SPM suggests that spatial heterogeneity greatly influences the drivers of Subjective Well-Being. This is particularly true for Life expectancy at birth, index of economic distress, illegal building rate and energy from renewable sources

4.3. Nowcasting of s80s20 index

In this paragraph we present the model for flash estimates implemented to provide updated data on the income inequality index, which is the ratio between the total income of the richest part of the population (first quintile: Q1) and the total income of the poorest one (fifth quintile: Q5).



	Panel Model						Spatial Panel Model								
	OLS			FE			RE			SAREM (ml)			SAREM (ml)		
	Coeff.	Std. Err.	Sig.	Coeff.	Std. Err.	Sig.	Coeff.	Std. Err.	Sig.	Coeff.	Std. Err.	Sig.	Coeff.	Std. Err.	Sig.
Intercept	70.080	95.608					-175.550	108.000		-124.320	73.594		-147.460	71.953	
Life expectancy at birth	-0.510	1.182		2.484	2.363		2.405	1.373		1.971	0.918	**	2.247	0.861	**
Early leavers from education and training	0.287	0.165		-0.393	0.180	*	-0.207	0.115	*	-0.208	0.123	*	-0.173	0.108	
Employment rate	0.555	0.113	***	-0.227	0.411		0.041	0.178	*	0.224	0.150	0.039	0.121		
Index of economic distress	-0.156	0.115		0.050	0.040		0.054	0.044		-0.001	0.055	0.032	0.057		
Voluntary activity	0.207	0.278		0.271	0.530		0.568	0.199		0.583	0.122	**	0.652	0.150	**
Rick-pocking rate	-0.449	0.158	***	-0.081	0.162		-0.171	0.120		-0.366	0.112	***	-0.167	0.123	
Impact of knowledge workers on employment	-1.335	0.327	***	-0.005	0.394		-0.091	0.319		-1.042	0.261	***	0.039	0.242	
Illegal building rate	-0.105	0.054	*	-0.007	0.071		-0.089	0.057		-0.133	0.042	**	-0.048	0.049	
Energy from renewable sources	-0.004	0.010		-0.014	0.018		0.028	0.007		0.016	0.009	*	0.051	0.008	
Irregularities in electric power distribution	2.430	0.815	***	0.617	0.491		0.621	0.471	***	1.226	0.582	**	0.525	0.539	
Year 2011															
Year 2012				2.617	0.773	***	2.915	0.712	***				2.919	0.708	***
Year 2013				-8.987	1.896	***	-8.568	1.391	***				-8.504	1.001	***
Year 2014				-10.478	2.687	***	-10.090	1.890	***				-9.859	1.275	***
Year 2015				-10.453	2.100	***	-9.897	1.565	***				-9.716	1.113	***
Year 2016				-11.671	2.676	***	-11.378	1.966	***				-11.172	1.260	***
Year 2017				-5.838	2.374	*	-5.420	1.687	***				-5.344	1.332	***

Figure 4.7: : Determinants of subjective well-being (first difference transformation in dependent variable)

Concerning the production process, the income inequality index is currently computed using the micro-data drawn from the European survey on *income and living conditions* (EU-SILC). Due to the data to be acquired, results are not available as timely as required by the policy evaluation cycle. For example at the end of 2018, the last available update was the 2016 one but, by the end of January 2019, Istat should provide the estimation of 2017 data¹.

However, this delay is common across European countries. Infact, providing timelier statistics on income poverty and inequality is a priority for the Commission and the European Statistical System.

To overcome this issue a new methodology based on microsimulation and macro-economic models has been put in place by Eurostat².

Istat has faced this new challenge adopting for the first round, referring to the preliminary estimation of 2015, a micro approach based on the microsimulation model developed by Istat (?). Starting from the estimation for 2016, Istat switched to a macro approach using as covariates the timelier information on the poverty rate and the saving rate.

For example, to estimate the 2017 data of the first quintile (Q1) the strategy has been to regress Q1 with the absolute poverty rate that is available 6 months after the end of the reference period.

Results are presented in the following equation, where $d15$ represents a dummy variable for the year 2015 that has been characterized by a very low level of Q1 compared to the average of previous years.

$$\widehat{q1_perc} = 7.79 - 0.369 d15 - 0.14 pov_ass$$

(0.099) (0.139) (0.0184)

$$T = 13 \quad \bar{R}^2 = 0.8799 \quad F(2, 10) = 44.947 \quad \hat{\sigma} = 0.1194$$

(standard errors in parentheses)

Figure 4.8: Q1 estimation based on poverty rate.Years 2007-2018

The estimation performance is reported in Figure 4.10

¹ And the end of March the estimation should be provided also for 2018

² See the website related to the experimental statistics <https://ec.europa.eu/eurostat/web/experimental-statistics/income-inequality-and-poverty-indicators>

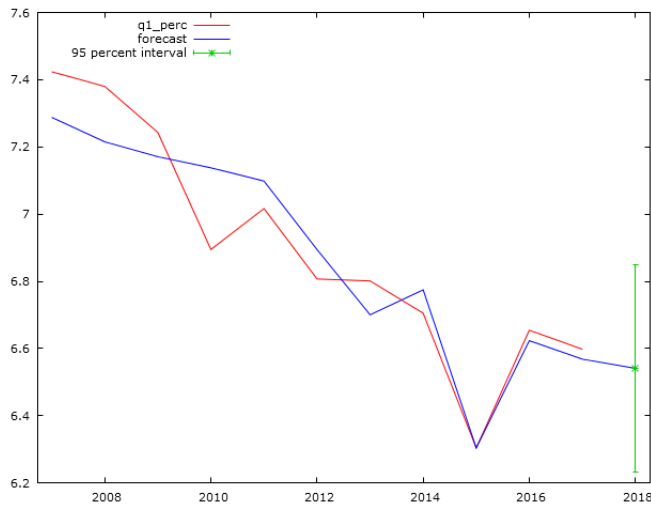


Figure 4.9: Q1 estimation based on poverty rate. Percentage points. Italy. Years 2007-2018

A similar approach has been followed for the estimation of Q5 using the saving rate (sav) drawn from the annual national accounts that are available 2 months after the reference year. In this case the equation started from 2007.

Here the results of the estimation:

$$\widehat{q5_perc} = 40.97 - \frac{0.18}{(0.031)} \text{ sav}$$

$T = 11 \quad \bar{R}^2 = 0.76 \quad F(1, 9) = 33.081 \quad \hat{\sigma} = 0.15637$
(standard errors in parentheses)

Figure 4.10: Q5 estimation based on saving rate. Italy. Years 2007-2018

while Figure 4.11 reports the graph for the real and the estimated values.

4.4. Daily adjustment for the social mood

Once defined the daily index some treatments are performed in order to highlight the properties useful for an analysis of the dynamic of the indicator. The daily frequency of the Social Mood on Economy Index implied that for the identification of the seasonal component was used a methodology different by the Istat standard, which is based on the model-based procedure implemented in TRAMO-SEATS. The treatment of seasonality in daily time series present, indeed, different methodological problems compared to the monthly or quarterly series due to the presence of multiple seasonality (midweek, weekly, monthly and/or annual) and the difficulty of distinguishing between the trend component and the annual seasonality. However, there are considerable advantages in analyzing daily time series with respect to their monthly aggregation, including greater availability of information and better identification of the effects of working days and other daily events that influence the dynamic of the series. Among the recent approaches adopted for the seasonal adjustment of high frequency series, in the analysis has been considered the methodology described in De Livera et al. (2011), that applies



Figure 4.11: Q5 estimation based on saving rate. Percentage points. Italy. Years 2007-2018

modified exponential smoothing models compared to the standard models introduced by Holt and Winters. The methodology is based on the hypothesis that each historical series can be represented as a combination of different components: a component that describes the level of the series, one for the growth rate (trend), one that captures seasonal movements and, finally, an irregular one. The model can be represented as:

$$\begin{aligned}
 y_t &= l_{t-1} + \phi b_{t-1} + \sum_{i=1}^T s_{t-m_i}^{(i)} + d_t \\
 l_t &= l_{t-1} + \phi b_{t-1} + \alpha d_t \\
 b_t &= \phi b_{t-1} + \beta d_t \\
 s_t^{(i)} &= s_{t-m_i}^{(i)} + \gamma_i d_t \\
 d_t &= \sum_{i=1}^p \psi_i d_{t-1} + \sum_{i=1}^q \theta_i e_{t-1}
 \end{aligned}$$

where m_1, \dots, m_T denote the seasonal periods, l_t is the local level in period t , b_t is the short-run trend in period t , $s_t^{(i)}$ represents the i -th seasonal component at time t , d_t denotes an ARMA(p, q) process, and e_t is a Gaussian white noise process with zero mean and constant variance σ^2 . The smoothing parameters are given by α , β , and γ_i for $i = 1, \dots, T$. The seasonal component is modeled through a trigonometric representation based on Fourier series. The ability to handle complex seasonality is a key advantage of the trigonometric approach over most traditional decomposition methods.

For the identification of these components, the model is then represented in a state space form whose parameters are estimated by the `tbats` function included in the `forecast` package in R. TBATS considers different alternatives and fit quite a few models: with Box-Cox transformation and without it, with and without trend and trend Damping, with and without ARMA(p, q) process used to model residuals and various amounts of harmonics are used to model the seasonal effects. The final model is chosen



using Akaike information criterion (AIC).

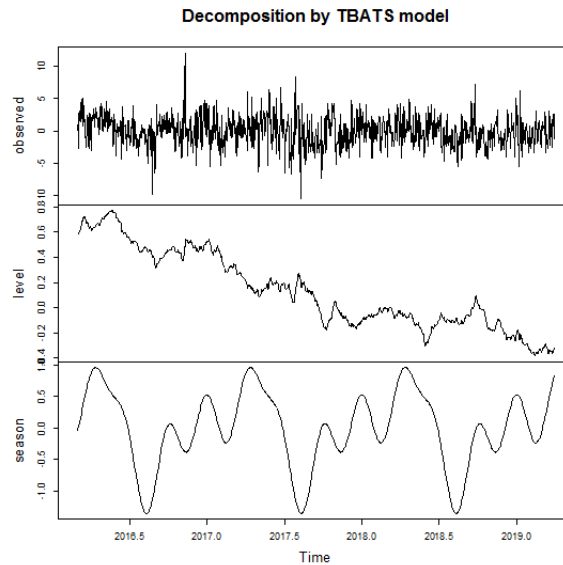


Figure 4.12: Decomposition by TBATS model

In particular, the seasonal adjustment of the daily time series of the Social Mood on Economy Index was implemented in two phases using the advantage of the framework proposed in Tbats that is able to encompass various deterministic effects that are often seen in real time series. The first step is the identification and estimation of deterministic effects through the introduction of appropriate dummy variables. In the second step the trend and seasonal component have been estimated using the linearized series. Figure 4.12 shows the result of seasonal adjustment procedure applied to the last release of social mood index 1/4/2016- 31/7/2019, which shows the linearized series of the index (observed) and the two components: the trend (level) and seasonal one (season). The series only shows a seasonality with an annual frequency, with positive peaks in January and in the period late March - early April and a negative peak in the summer months.



5. Discussion

Multivariate methods and Dutch Labour Force

Since 2010 official monthly figures of the Dutch Labour Force are produced with a multivariate structural time series model. The sample size of the Dutch LFS is too small to use a direct estimator like the GREG estimator for the production of monthly figures since the variance of these estimates are unacceptable large. With the time series model, sample information from preceding sample occasions is used to produce more accurate model-based estimates. The model also accounts for the rotation group bias and autocorrelation induced by the rotating panel design of the Dutch LFS.

To further improve the accuracy as well as the timeliness of the monthly estimates this model is extended with available auxiliary information. Potential auxiliary series in the Netherlands are claimant counts and Google Trends. Claimant counts is an univariate series that can be added as an additional series to the time series model of the LFS in straightforward manner. The time series estimates of LFS are improved with the claimant counts series by modelling the correlation between the slope disturbance terms of the trend component of the LFS and the claimant counts. With the claimant counts series the precision of the monthly LFS estimates is improved with about 12%. The claimant counts for period t become available in period $t + 1$. Therefore it was anticipated that claimant counts are not a potential auxiliary series for nowcasting. Nevertheless the one-step ahead forecasts for the LFS estimates are about 25% more accurate compared to a time series model that only uses LFS time series.

A second source of auxiliary series are Google Trends for search terms that are expected to be related with unemployment. Google Trends can be downloaded at a higher frequency, e.g. on a weekly basis. This results in more timely related information, which can be used to make nowcasts for the LFS parameters in period t at the moment that the Google Trends become available in this period, while the survey data are not available yet. In addition these data sources can be used to further improve the time series model estimates for the LFS in a similar as with the claimant counts.

An issue with a big data source like Google Trends is that easily a large amount of auxiliary series can be derived. Including them in a standard multivariate structural time series model that accounts for the correlations between e.g. the trend disturbance terms result in a high dimensionality problem. Therefore a dynamic factor model is proposed, which allows to combine a large set of auxiliary series with the target series from the survey sample in a parsimonious model. The approach resembles the method originally proposed by Doz et al. (2011). The method is based on a two-step estimator. In a first step a small number of common factors are derived from the Google Trends using principal component analysis. In a second step the Google Trends are combined with the LFS series in structural time series model, where the common factors of the Google Trends are modelled as a local level model and re-estimated with the Kalman filter. The factor loadings obtained in the first step are kept fixed in the design matrix of the measurement equation. The estimates for the LFS are improved by modelling the correlation between slope disturbance terms of the trend for the LFS with the level disturbance



terms of the Google Trend common factors.

The rotating panel design of the LFS results in a rather complex state-space model for the sampling error. This hampers the application of a model that temporarily disaggregates the LFS to a weekly frequency. Therefore we were forced to apply a model that temporarily aggregates the Google Trends to a monthly frequency. In other applications, where the model for the survey time series is less complex, temporal disaggregation of the survey time series is recommended.

Results from a likelihood ratio test are in favour of a model that contains claimant counts rather than Google Trends. Nonetheless, the accuracy of the estimation and of the nowcast of the level and the change in unemployment, does not deteriorate when only the latter series are included. A model that combines the LFS series with the Google Trends only slightly improves the out-of-sample forecasts. The out-of-sample forecasts are, however, more accurate if the LFS series are combined with claimant counts, despite the fact that for both series the observation for period t becomes available in period $t + 1$.

The results with the Google Trends slightly improve when the cycle of the factor is appropriately modelled according to a (seasonal) ARIMA model. The change in unemployment does not seem to depend on the lagged Google Trends, suggesting that job-related terms are not searched on Google before becoming unemployed. Targeting the search terms with the Elastic Net before estimating their factors does not improve the estimation and nowcast accuracy of the unobserved components.

Our choice of search terms is hand-picked, and therefore to some extent arbitrary and limited. We can therefore not rule out the usefulness of Google Trends for unemployment estimation in the Netherlands, although our results suggest limited scope for improvement. Clearly, for other topics or other countries results might be entirely different. A first observation is that for countries with strongly related series from a register like claimant counts, the Google Trends series do not add much additional information for nowcasting unemployment. For countries without such register information, it might still be worthwhile to explore the possibilities of Google Trends.

Further research is focused on time series models that allow for time varying correlations between the disturbance terms of the stochastic trend and seasonal components. For claimant counts, changes in the legislation who is qualified to receive unemployment benefits might result in gradual changing correlations between the LFS series and the claimant counts series. For Google Trends search behaviour on job related search terms before the start of the financial crisis in 2009 might be completely different compared to the period directly after the financial crisis, which might disturb the relation with the LFS series over time. It might be expected that applying a time series model that incorrectly assumes a time invariant correlation to a set of series where the correlation is changing over time, results in biased parameter estimates (van den Brakel and Krieg, 2016). Statistics Netherlands currently investigates different methods to relax the assumption that correlation are time invariant, e.g. by using Generalized Autoregressive Score models (Creal et al., 2013, Harvey, 2013).



Dynamic factor model, space panel model and well-being

Till now analysis and methodology has been more concentrated on cross-section approach. Despite the importance of spatial and temporal dimensions, there is a lack of study in the relationship between SWB and its determinants;

The availability over time of Well-Being data over time should be increased by NSOs or other statistical agency.

Dynamic factor models have been used to explore their potentiality on the analysis of well-being. Results from DFM and SPM highlight a high degree of spatial heterogeneity in well-being across Italian regions:

DFM shows that the main factor captures different domains of well-being according to the territorial levels of analysis: Satisfaction with family relations and friendship are the most important factor across Italy; Satisfaction with means of transport, obesity and adequate nutrition are important mainly for Piedmont region, Sedentarioussness and indicators related to the economic well-being are most important in the Sicily region (in the South of Italy)

SPM shows that spatial heterogeneity greatly influences the drivers of Subjective Well-Being. This is particularly true for Life expectancy at birth, index of economic distress, illegal building rate and energy from renewable sources.

More general our analysis suggest that time series approach could be applied, thus allowing to shed light both on heterogeneity across regions and on the interaction across the indicators;

However switching to a time series approach requires the traditional controls on the variables such as log transformation and integration.



Bibliography

- Alonso, A. M., J. Rodríguez, C. García-Martos, and M. Jesús Sánchez (2011). Seasonal dynamic factor analysis and bootstrap inference: Application to electricity market forecasting. *Technometrics* 53, 137–151.
- Askatas, N. and K. Zimmermann (2009). Google econometrics and unemployment forecasting. *Applied Economics Quarterly* 55, 107–120.
- Bacchini, F., B. Baldazzi, R. De Carli, L. Di Biagio, M. Savioli, M. P. Sorvillo, and A. Tinto (2018). The Italian framework to measure well-being: towards the 2.0 version. *Submitted*.
- Bai, J. (2004). Estimating cross-section common stochastic trends in nonstationary panel data. *Journal of Econometrics* 122, 137–183.
- Bai, J. and S. Ng (2008). Forecasting economic time series using targeted predictors. *Journal of Econometrics* 146, 304–317.
- Bailar, B. (1975). The effects of rotation group bias on estimates from panel surveys. *Journal of the American Statistical Association* 70, 23–30.
- Barigozzi, M. and M. Luciani (2017). Common Factors, Trends, and Cycles in Large Datasets. Finance and economics discussion series 2017-111.
- Boivin, J. and S. Ng (2005). Understanding and comparing factor-based forecasts. *International Journal of Central Banking* 3, 117–151.
- Choi, H. and H. Varian (2009). Predicting initial claims for unemployment benefits. *Google Inc*, 1 – 5.
- Creal, D., S. Koopman, and A. Lucas (2013). Generalized autoregressive score models with applications. *Journal of Applied Econometrics* 28, 777 – 795.
- Datta, G., P. Lahiri, T. Maiti, and K. Lu (1999). Hierarchical bayes estimation of unemployment rates for the states of the u.s. *Journal of the American Statistical Association* 94(448), 1074–1082.
- De Livera, A. M., R. J. Hyndman, and R. D. Snyder (2011). Forecasting time series with complex seasonal patterns using exponential smoothing. *Journal of the American Statistical Association* 106(496), 1513–1527.
- Doornik, J. (2009). *An Object-oriented Matrix Programming Language Ox 6*. Timberlake Consultants Press: London.
- Doz, C., D. Giannone, and L. Reichlin (2011). A two-step estimator for large approximate dynamic factor models based on kalman filtering. *Journal of Econometrics* 164, 188–205.
- Durbin, J. and S. Koopman (2012). *Time Series Analysis by State Space Methods*. Oxford University Press.



- Durbin, J. and B. Quenneville (1997). Benchmarking by state space models. *International Statistical Review* 65, 23–48.
- Elliott, G., T. J. Rothenberg, and J. H. Stock (1996). Efficient tests for an autoregressive unit root. *Econometrica* 64(4), 813–836.
- Fay, R. and R. Herriot (1979). Estimates of income for small places: An application of james-stein procedures to census data. *Journal of the American Statistical Association* 74(366), 269–277.
- Forni, M., A. Giovannelli, M. Lippi, and S. Soccorsi (2018). Dynamic factor model with infinite-dimensional factor space: Forecasting. *Journal of Applied Econometrics* 33(5), 625–642.
- Giannone, D., L. Reichlin, and D. Small (2008). Nowcasting: The real-time informational content of macroeconomic data. *Journal of Monetary Economics* 55, 665–676.
- Harvey, A. (1989). *Forecasting, structural time series models and the Kalman filter*. Cambridge University Press.
- Harvey, A. (2013). *Dynamic models for volatility and heavy tails: with application to financial and econometric time series*. Cambridge University Press.
- Harvey, A. and C. Chung (2000). Estimating the underlying change in unemployment in the uk. *Journal of the Royal Statistitcal Society, A series* 163, 303–339.
- Hastie, T. and H. Zou (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67, 301–320.
- Koopman, S., A. Harvey, N. Shephard, and J. Doornik (2009). *STAMP 8.2; Structural Time Series Analyser, Modeller and Predictor*. Timberlake Consultants Press: London.
- Koopman, S., N. Shephard, and J. Doornik (2008). *SsfPack 3.0: Statistical Algorithms for Models in State Space Form*. Timberlake Consultants Press: London.
- Koopman, S., N. Shephard, and J. Doornik (2009). Statistical algorithms for models in state space form using ssfpack 2.2. *Econometrics Journal* 2, 113–166.
- Krieg, S. and J. van den Brakel (2012). Estimation of the monthly unemployment rate for six domains through structural time series modelling with cointegrated trends. *Computational Statistics and Data Analysis* 56, 2918–2933.
- Lazer, D., R. Kennedy, G. King, and A. Vespignani (2014). The parable of Google Flu: traps in big data analysis. *Science* 343, 1203–1205.
- Levin, A., C.-F. Lin, and C.-S. J. Chu (2002). Unit root tests in panel data: asymptotic and finite-sample properties. *Journal of econometrics* 108(1), 1–24.
- Marcellino, M., J. Stock, and M. Watson (2003). Macroeconomic forecasting in the euro area; country specific versus area-wide information. *European economic review* 47, 1–18.
- Moauero, F. and G. Savio (2005). Temporal disaggregation using multivariate structural time series models. *Econometrics Journal* 8, 214–234.



- Moon, H. R. and B. Perron (2012). Beyond panel unit root tests: Using multiple testing to determine the nonstationarity properties of individual series in a panel. *Journal of Econometrics* 169, 29–33.
- Pfeffermann, D. (1991). Estimation and seasonal adjustment of population means using data from repeated surveys. *Journal of Business & Economic Statistics* 9, 163–175.
- Pfeffermann, D. and L. Burck (1990). Robust small area estimation combining time series and cross-sectional data. *Survey Methodology* 16, 217–237.
- Pfeffermann, D. and R. Tiller (2006). Small area estimation with state-space models subject to benchmark constraints. *Journal of the American Statistical Association* 101, 1387–1397.
- Rao, J. and M. Yu (1994). Small area estimation by combining time series and cross-sectional data. *The Canadian Journal of Statistics* 22, 511–528.
- Schiavoni, C., F. Palm, S. Smeekens, and J. van den Brakel (2019). A dynamic factor model approach to incorporate big data in state space models for official statistics. Discussion paper statistics netherlands and maastricht university.
- Stephens-Davidowitz, S. and H. Varian (2015). A hands-on guide to google data. *Google, Inc.*, 1 – 25.
- Stiglitz, J., A. Sen, J.-P. Fitoussi, et al. (2009). The measurement of economic performance and social progress revisited. *Reflections and overview. Commission on the Measurement of Economic Performance and Social Progress, Paris.*
- Stock, J. and M. Watson (2002a). Forecasting using principal components from a large number of predictors. *Journal of the American Statistician Society* 97, 1167–1179.
- Stock, J. and M. Watson (2002b). Macroeconomic forecasting using diffusion indexes. *Journal of Business and Economic Statistics* 20, 147–162.
- Team, R. C. (2017). R: A language and environment for statistical computing. R foundation for statistical computing, vienna, austria.
- Tinto, A., F. Bacchini, B. Baldazzi, A. Ferruzza, J. Van den Brakel, R. Willems, N. Rosenski, T. Zimmermann, Z. Andrási, M. Farkas, and Z. Fábián (2018). Report on international and national experiences and main insight for policy use of well-being and sustainability framework. Technical report, MAKSWELL—deliverable 1.1. https://www.makswell.eu/attached_documents/output_deliverables/deliverable_1.1_draft.pdf.
- van den Brakel, J. and S. Krieg (2015). Dealing with small sample sizes, rotation group bias and discontinuities in a rotating panel design. *Survey Methodology* 41, 267–296.
- van den Brakel, J. and S. Krieg (2016). Small area estimation with state space common factor models for rotating panels. *Journal of the Royal Statistical Society, A series* 179, 763–791.
- van den Brakel, J., P. Smith, N. Tzavidis, R. Iannaccone, D. Zurlo, F. Bacchini, L. Di Consiglio, T. Tuoto, , M. Pratesi, C. Giusti, S. Marchetti, S. Bastianoni, G. Betti, A. Lemmi, F. Pulselli, and L. Neri (2019). Methodological aspects of using big-data, makswell, wp2, deliverable 2.2.



- van den Brakel, J., E. Söhler, P. Daas, and B. Buelens (2017). Social media as a data source for official statistics; the dutch consumer confidence index. *Survey Methodology* 43(2), 183–210.
- You, Y. (2008). An integrated modeling approach to unemployment rate estimation for sub-provincial areas of canada. *Survey Methodology* 34(1), 19–27.
- You, Y., J. Rao, and J. Gambino (2003). Model-based unemployment rate estimation for the canadian labour force survey: A hierarchical bayes approach. *Survey Methodology* 29(1), 25–32.
- Zardetto, D. (2018). Using twitter data for the social mood on economy index. 13 conferenza nazionale di statistica.