



**HAL**  
open science

# Regression analysis with missing data and unknown colored noise: application to the MICROSCOPE space mission

Quentin Baghi, G. Métris, Joël Bergé, Bruno Christophe, Pierre Touboul,  
Manuel Rodrigues

## ► To cite this version:

Quentin Baghi, G. Métris, Joël Bergé, Bruno Christophe, Pierre Touboul, et al.. Regression analysis with missing data and unknown colored noise: application to the MICROSCOPE space mission. *Physical Review D*, 2015, 91, pp.062003. 10.1103/PhysRevD.91.062003. hal-01111300v2

**HAL Id: hal-01111300**

**<https://hal.science/hal-01111300v2>**

Submitted on 12 Mar 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Regression analysis with missing data and unknown colored noise: application to the MICROSCOPE space mission

Quentin Baghi,<sup>1,\*</sup> Gilles Métris,<sup>2,†</sup> Joël Bergé,<sup>1</sup> Bruno Christophe,<sup>1</sup> Pierre Touboul,<sup>1</sup> and Manuel Rodrigues<sup>1</sup>

<sup>1</sup>ONERA - The French Aerospace Lab, 29 avenue de la Division Leclerc, 92320 Chatillon, France

<sup>2</sup>Geoazur (UMR 7329), Observatoire de la Côte d’Azur Bt 4,

250 rue Albert Einstein, Les Lucioles 1, Sophia Antipolis, 06560 Valbonne, France

(Dated: February 17, 2015)

The analysis of physical measurements often copes with highly correlated noises and interruptions caused by outliers, saturation events or transmission losses. We assess the impact of missing data on the performance of linear regression analysis involving the fit of modeled or measured time series. We show that data gaps can significantly alter the precision of the regression parameter estimation in the presence of colored noise, due to the frequency leakage of the noise power. We present a regression method which cancels this effect and estimates the parameters of interest with a precision comparable to the complete data case, even if the noise power spectral density (PSD) is not known *a priori*. The method is based on an autoregressive (AR) fit of the noise, which allows us to build an approximate generalized least squares estimator approaching the minimal variance bound. The method, which can be applied to any similar data processing, is tested on simulated measurements of the MICROSCOPE space mission, whose goal is to test the Weak Equivalence Principle (WEP) with a precision of  $10^{-15}$ . In this particular context the signal of interest is the WEP violation signal expected to be found around a well defined frequency. We test our method with different gap patterns and noise of known PSD and find that the results agree with the mission requirements, decreasing the uncertainty by a factor 60 with respect to ordinary least squares methods. We show that it also provides a test of significance to assess the uncertainty of the measurement.

**PACS numbers:** 04.80.Cc, 04.80.Nn, 07.87.+v, 95.55.-n, 07.05.Kf

**Keywords:** Data processing, Experimental test of gravitational theories, Spaceborne instruments

## I. INTRODUCTION

Situations where series of measurements, ideally regularly sampled, suffer from short interruptions are common in a wide range of applications and experimental set-ups. It is also usual to perform linear regression analysis of data samples, in order to estimate parameters of interest by fitting other data series to the measured signals. In particular, this is a typical scenario for space missions in fundamental physics such as MICROSCOPE [1, 2] and LISA Pathfinder [3]. Long time integrations are needed by these experiments to reach the required signal-to-noise ratios (SNR) or the required levels of free-fall at the frequencies of interest. The duration of such measurements increases the probability to have invalid data in the integration period. It has been found that gaps could arise in the time series measured by the accelerometers carried on-board the MICROSCOPE satellite, and that those gaps could have substantial impact on the outcome of the regression when data is noisy.

Here “gaps” refers to either lack of data or unusable information such as saturations and outliers during short or long time spans, which are eventually discarded. In the case of the MICROSCOPE space mission, discontinuities in the data availability could be due to data losses in the telemetry transmission, while data alteration could be

the consequence of three main identified causes: crackles in the cold gas tanks triggered by decreasing pressure as they empty, crackles in the multi-layer insulation (MLI) coating due to temperature variations in flight, or micro-meteorites impacts. All saturated data are clearly identified by a flagging system in the telemetry.

The objective of the MICROSCOPE signal processing can be regarded as rather general. It consists in detecting and estimating the amplitude of a periodic signal present in some measured time series. In the studied case the signal is the signature of a possible violation of the Weak Equivalence Principle (WEP), as detailed later, and is expected to arise around a certain frequency that we denote  $f_{EP}$ . The amplitude to be estimated is the “EP parameter”, denoted  $\delta$ . In previous works [4] the data analysis had been optimized in order to minimize the projection of possible unknown harmonic perturbations onto the signal of interest by an appropriate tuning of its frequency  $f_{EP}$  and/or the integration duration, in particular in the case of missing data. At the time, instrumental noise had been disregarded in order to exclusively deal with projection effects. Here we rather focus on the impact of missing data on the noise affecting the estimation.

While the proposed approach is applied to MICROSCOPE simulated data, it leads to provide a robust method to estimate one or several deterministic components in the general context of time series with missing data affected by unknown colored noise. Although we have physical models of the expected noise spectrum, we assume in this study that it is not known *a priori*, allow-

---

\* quentin.baghi@onera.fr

† gilles.metris@oca.eu

ing us to cope with the most general situation.

We show that noise distortions due to missing data points may dramatically increase the uncertainty of the estimation. This is due to the convolution effect between the observation window and the original noise spectrum, which leads to a leakage of the frequencies where the power is high to the frequencies where the power is low.

Methods such as Ordinary Least Squares (OLS) or equivalently Lomb-Scargle periodogram [5, 6], as well as CLEAN-like algorithms [7], may fail in retrieving the required precision [8, 9], mainly because these approaches rely on a white noise assumption. In order to increase the precision of the fit, the noise correlation matrix must also be estimated. A general approach is to maximize the likelihood function with respect to both regression parameters ( $\delta$  in our problem) and the noise correlation matrix. Such an approach can use the Expectation-Minimization (EM) procedure like MAPES algorithms [10]. However, their convergence may be very slow, especially for large data samples like in the MICROSCOPE case (about  $10^6$  points). More recent works also use least squares iterative adaptive approaches (IAA) to estimate harmonic and noise parameters iteratively [11], but require to store and invert correlation matrices, which is computationally expensive with an observation vector of  $10^6$  entries. Likewise, the authors of the last two techniques do not present applications with colored noise. Some methods are already implemented to extract unknown colored spectral densities, especially in the domain of gravitational waves detection (see for example [3, 12–14]), but they do not tackle the problem of gapped time series. A suitable method is thus developed to estimate the EP parameter in case of missing data.

Another type of algorithms referred to as “inpainting” techniques is based on a sparsity-prior to fill the gaps [15, 16]. Their adaptation to general noise spectra is currently studied in the MICROSCOPE team (Bergé *et al.*, in prep.). We rather focus here on an approach that avoids filling the gaps.

We develop a method with two successive objectives. The first one is to reach the order of magnitude of the original (i.e. complete data) uncertainty in the estimation of the amplitudes of the deterministic components we are looking for. The second objective is to theoretically quantify the improvement on the variance of the estimator, using an approach that does not require to fill in the data gaps.

Our technique is based on the estimation of the noise spectrum by using a high-order autoregressive (AR) model. The result is used to weight the data through an orthogonalization of the covariance matrix. This leads to an approximation of the best estimator in the sense of the variance, also referred to as the Best Linear Unbiased Estimator (BLUE) which is also the Generalized Least Squares (GLS) estimator in a linear regression context. The main idea in the proposed approach is to separately estimate the noise coefficients and the regression parameters instead of jointly estimating all the parameters. This

is done in an iterative procedure that avoids the use of non-linear optimization algorithms.

The proposed approach, that we dub “Kalman-AR Model Analysis” or “KARMA” for short, is divided in three steps. The first step consists in estimating the AR parameters describing the noise. This is done by using Burg’s algorithm adapted to discontinuous data [17]. The second step is carried out via a Kalman filter algorithm based on the AR model that allows us to compute the weights, as shown by Jones [18]. In the third step we finally compute an approximation of the Generalized Least Squares (GLS) estimator of the regression parameters, in a way similar to maximum likelihood computation methods applied to regression models [19, 20]. These steps can be reproduced to converge to the maximum likelihood estimator (MLE) of the parameters.

In this paper, we first analyze the effect of the missing data pattern on the estimation uncertainty (section II). We then describe the KARMA method (section III) and we present a way of evaluating its performance, allowing us to give a criteria for the detection of the searched signal (section IV). Finally, after a brief description of the mission context, we apply this technique to MICROSCOPE simulated time series, in particular to data samples generated with the mission and instrument simulator (section V). In section VI we discuss the results.

## II. IMPACT OF MISSING DATA

Although we apply our study to the MICROSCOPE data analysis, it can be viewed as a general regression problem. The measurement equation can be summarized as follows:

$$\boldsymbol{\gamma} = \delta \mathbf{s}_{\text{EP}} + \sum_i \alpha_i \mathbf{s}_{p,i} + \mathbf{z}, \quad (1)$$

where  $\boldsymbol{\gamma}$  is the  $N$ -points complete measurement vector defined as  $\boldsymbol{\gamma} = (\gamma_0 \dots \gamma_{N-1})^T$ , and  $\delta$  and  $\mathbf{s}_{\text{EP}}$  are respectively the parameter and the signal of interest (the EP parameter and the EP violation signal for our purpose).

The second term accounts for possible perturbations, whose amplitudes  $\alpha_i$  should also be estimated to reject any bias.

The third term is the residual noise vector  $\mathbf{z}$  assumed to be a zero-mean Gaussian random vector. The main objective is the estimation of  $\delta$ , for which the square root of the one-sided noise power spectral density (PSD) at EP frequency must be  $1.4 \times 10^{-12} \text{ ms}^{-2}/\sqrt{\text{Hz}}$  [1].

The presence of missing or corrupted data in the time series is identified by a mask vector  $\mathbf{w}$  which is equal to 1 when the data is available and 0 otherwise, regardless of the nature of the gap. The observed signal is thus the vector  $\mathbf{y}$  with entries  $y_n = w_n \gamma_n$ . We assume that the loss of data arises before any possible filtering.

### A. Impact on the PSD

We briefly derive the impact of the observation window  $w$  on the PSD of a pure stationary random signal. Thus in this section we assume  $\gamma = z$ . The real signal  $\gamma$  is regularly sampled at a frequency  $f_s$  so that  $\gamma_n = \gamma(n/f_s)$ .

For a stationary discrete parameter process, the autocovariance function is defined as:

$$R_y(k) \equiv \text{E}[y_n y_{n+k}] - \text{E}[y_n] \text{E}[y_{n+k}]. \quad (2)$$

Then the PSD is the Discrete-Time Fourier Transform (DTFT) of the autocovariance [21]:

$$S_y(f) = \frac{1}{f_s} \sum_{k=-\infty}^{+\infty} R_y(k) e^{-2j\pi k f / f_s}. \quad (3)$$

In the case of the masked noise  $y_n = w_n z_n$ , Eq. (2) gives:

$$R_y(k) = \text{E}[w_n z_n w_{n+k} z_{n+k}] - \text{E}[w_n z_n] \text{E}[w_{n+k} z_{n+k}].$$

We assume that the underlying process in  $z$  is independent of the window  $w$ , and that  $w$  is a stationary process, so that one can write:

$$\begin{aligned} R_y(k) &= \text{E}[z_n z_{n+k}] \text{E}[w_n w_{n+k}] - \mu_z^2 \mu_w^2 \\ &= (R_z(k) + \mu_z^2) (R_w(k) + \mu_w^2) - \mu_z^2 \mu_w^2, \end{aligned} \quad (4)$$

where for any random variables  $x$  we note  $\mu_x$  its expectation.

Assuming that  $z_n$  is a zero-mean process ( $\mu_z = 0$ ), the PSD of the windowed signal is obtained by taking the DFT of Eq. (4):

$$S_y(f) = \mu_w^2 S_z(f) + [S_w * S_z](f), \quad (5)$$

where  $*$  is the convolution operator.

The first term can be viewed as a loss of power due to the missing data and the second term accounts for the frequency leakage. In the case of uniform random gaps, one shows (see appendix A) that  $\mu_w$  is equal to the probability to have a gap at a given time. Then  $S_w(f)$  is a constant, and the leakage term is proportional to the mean power. Therefore, the noise will increase significantly at frequencies where the leakage term is dominant.

As an illustration, a simulation of the MICROSCOPE instrumental noise alone is presented in Fig. 1. The noise is generated using an approximate PSD model, taking into account thermal sensitivities at lower frequencies, position sensor noise at higher frequencies, random noise of the pick-up circuitry and the frequency response of the control loop:

$$2S_z(f) = \sigma_z^2 \left( 1 + \left( \frac{f}{f_1} \right)^{-1} + \left( \frac{f}{f_2} \right)^4 \right) \cdot |H_{cl}(f)|^2 \quad (6)$$

with  $\sigma_z = 1.4 \times 10^{-13} \text{ ms}^{-2}/\sqrt{\text{Hz}}$ ,  $f_1 = 8.1 \times 10^{-2} \text{ Hz}$  and  $f_2 = 1.3 \times 10^{-2} \text{ Hz}$ .  $H_{cl}$  is the transfer function

of the closed control loop of the accelerometer. It has almost a unit gain for all frequencies under 1 Hz, and induces a slight inflection in higher frequencies. The factor of 2 accounts for the fact that  $S_z(f)$  is the two-sided PSD. The data is sampled at a frequency  $f_s = 4 \text{ Hz}$  on a duration  $T = 1.4 \text{ days}$  corresponding to 20 satellite orbits with  $N_g = 5200$  gaps of the same length (0.5 seconds), randomly distributed over the time series. We observe a transfer of power from high frequencies to low frequencies, increasing the apparent noise around  $f_{\text{EP}} = 9.4 \times 10^{-4} \text{ Hz}$  by two orders of magnitude.

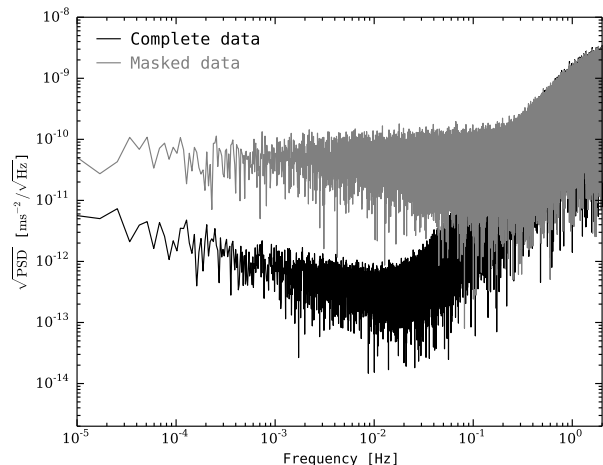


FIG. 1. Periodogram of original (black) and incomplete (grey) time series with 0.5 second data gaps randomly distributed in a 20 orbits session. The simulation is done for 260 random gaps per orbit.

### B. Impact on the least squares estimate

We now demonstrate that the observed increase of the noise is not a simple artifact of the Fourier representation but directly impacts the estimation uncertainty in a least squares fitting approach. We assume that the analyzed signal is the sum of a harmonic component  $s_{\text{EP}}$  at frequency  $f_{\text{EP}}$  and a correlated Gaussian random noise  $z$ . For the sake of simplicity, we ignore the presence of possible deterministic perturbations, therefore the signals  $s_{p,i}$ 's in Eq. (1) are all zero. The signal is still sampled at frequency  $f_s$  on  $N$  data points. Thus the signal reads:

$$\gamma = \delta s_{\text{EP}} + z. \quad (7)$$

We define the window matrix as the diagonal matrix formed by the window vector:  $W = \text{diag}(w_0 \dots w_{N-1})$ . We aim at calculating the variance of the OLS estimate that only uses the available data (at times for which  $w_n = 1$ ). In the least squares formalism, this is equivalent to studying the windowed vector  $y = W\gamma$ . We also define the model matrix  $A$ . Although

it can take a general form including various signals, we assume here that it contains the EP signal model only such that  $A = \mathbf{s}_{\text{EP}}$ . We also define the masked model matrix  $A_w = WA$ . The usual OLS formulas give the following parameter estimate:

$$\hat{\delta} = (A_w^\dagger A_w)^{-1} \cdot A_w^\dagger \mathbf{y}, \quad (8)$$

as well as its variance:

$$\text{Var}(\hat{\delta}) = K^{-1} A_w^\dagger \Sigma_y A_w K^{-1}, \quad (9)$$

where we defined  $K = A_w^\dagger A_w$  and  $\Sigma_y = W \Sigma_z W^\dagger$  with  $\Sigma_z = \text{E}[\mathbf{z} \mathbf{z}^\dagger]$ , the covariance matrix of the noise vector. Here  $\dagger$  denotes the hermitian adjoint. As a result, the noise correlation seen by the estimator is  $\Sigma_y$  instead of  $\Sigma_z$  in the complete case.

In the case of a stationary Gaussian random noise the estimator covariance can be diagonalized in the Fourier space:

$$\Sigma_y = \frac{f_s}{N} M^\dagger D M, \quad (10)$$

where  $D$  is the diagonal matrix formed by the two-sided discrete PSD:  $D = \text{diag}(\hat{S}_0 \dots \hat{S}_{N-1})^T$  and  $M$  is the Discrete Fourier Transform (DFT) matrix with coefficients:  $M_{kl} = \exp(-\frac{2i\pi kl}{N})$ . The discrete spectrum is defined as the expectation of the periodogram. It can be seen as an approximation of the real PSD [21]:

$$\hat{S}_k \equiv \frac{1}{f_s} \sum_{n=-(N-1)}^{N-1} \left(1 - \frac{|n|}{N}\right) R_y(n) e^{-2j\pi \frac{nk}{N}}. \quad (11)$$

This diagonalization thus links the estimator variance and the PSD of the windowed noise. By developing Eq. (9) we show (see appendix B for more details) that in the case of a harmonic model such as  $\mathbf{s}_{\text{EP},n} = \gamma_{\text{EP}} \sin(2\pi n f_{\text{EP}}/f_s)$ , for sufficiently large  $N$ , the estimator variance is approximately equal to:

$$\text{Var}(\hat{\delta}) \approx \frac{2f_s N S_y(f_{\text{EP}})}{N_o^2 \gamma_{\text{EP}}^2}, \quad (12)$$

where  $S_y$  is given by equation (5),  $N_o = N - N_g$  is the number of observed data and  $\gamma_{\text{EP}}$  is the amplitude of the model, which is the gravitational acceleration in our case. As a result, in presence of missing data, the estimation variance increases proportionally to the leakage term in Eq. (5). To quantify the increase of the uncertainty, we plot the standard deviation of the estimator as a function of the number of data gaps per orbit in Fig. 2, in the case of short random gaps of fixed length (0.5 second) uniformly distributed over the time series (the effect of the size and the number of gaps is discussed in Bergé *et al.*, in prep.). The theoretical standard deviation (black curve) is obtained using Eq. (9). In order to check the correctness of the distribution, we also plot the sample standard deviation of 400 estimates (red curve)

corresponding to different realizations of the noise vector  $\mathbf{z}$ . This shows that the uncertainty grows by one order of magnitude from 10 gaps per orbit only, which represents a data loss of 0.04 %. This is not acceptable with respect to the performance objectives of the mission. Therefore an alternative estimation method needs to be implemented.

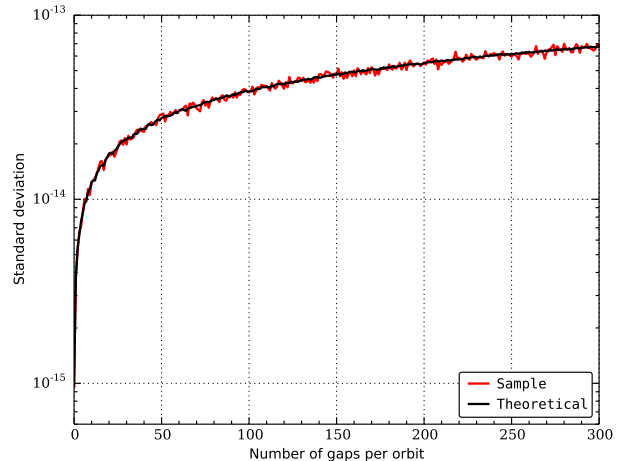


FIG. 2. Theoretical (black) and sample (red) standard deviations of the original least squares estimate of the EP parameter as a function of the number of gaps per orbits. All gaps have the same duration of 0.5 second and are randomly distributed over a 20 orbits session.

### III. KALMAN-AR MODEL ANALYSIS (KARMA)

The poor performance of the OLS estimator is due to the fact that its variance is not minimal. To minimize the variance, the Best Linear Unbiased Estimator (BLUE) is needed, which takes the form of a Generalized Least Squares (GLS) estimator in linear regression problems. In case of missing data, it reads:

$$\hat{\beta} = (A_o^\dagger \Sigma_o^{-1} A_o)^{-1} \cdot A_o^\dagger \Sigma_o^{-1} \mathbf{y}_o, \quad (13)$$

where  $\beta$  is the  $q \times 1$  vector of parameters to be estimated. The observation vector  $\mathbf{y}_o \equiv (\gamma_{n_0} \dots \gamma_{n_{N_o-1}})^T$  gathers the available data only, that is,  $n_0, \dots, n_{N_o-1}$  are the time indexes corresponding to the observed data. Similarly,  $A_o$  is the model matrix where we have kept only rows corresponding to observed data. Here  $A_o$  is assumed to be general, of size  $N \times q$ .  $\Sigma_o$  is the covariance matrix of the observed noise vector  $\mathbf{z}_o$  and admits a Cholesky decomposition such that  $\Sigma_o = L_o L_o^\dagger$  where  $L_o$  is a lower triangular matrix.

The difficulty here is to estimate the noise covariance matrix  $\Sigma_o$  in spite of the missing data. The method that we propose consists in calculating an approximation

of the GLS estimator by postulating an autoregressive (AR) model for the noise. This is done in three steps which are detailed below: estimation of AR parameters (step 1), calculation of the whitened vectors  $L_o^{-1}\mathbf{y}_o$  and  $L_o^{-1}A_o$  (step 2), calculation of the estimate  $\hat{\boldsymbol{\beta}}$  (step 3). The process may be iterated if necessary.

### A. Step 1: AR parameters estimation

The first step is to estimate the noise characteristics encapsulated in the covariance matrix  $\Sigma_z$ . To do so, we assume that the noise process can be described by an autoregressive model of some order  $p$  to be determined, verifying the following relation at all times  $n$ :

$$z_n + a_1 z_{n-1} + \dots + a_p z_{n-p} = \epsilon_n, \quad (14)$$

where  $a_1, \dots, a_p$  are the AR coefficients and  $\epsilon$  is a zero-mean white Gaussian random field of variance  $\sigma^2$ . Note that this is equivalent to approximating the noise PSD with a rational function, the numerator being a polynomial of degree  $p$  in  $\exp(-2i\pi f/f_s)$  such as:

$$\hat{S}_z(f) = \frac{\sigma^2/f_s}{|1 + a_1 e^{-2i\pi f/f_s} + \dots + a_p e^{-2i\pi p f/f_s}|^2}. \quad (15)$$

The choice of this model is motivated by the following arguments. Firstly the use of a parametric model consistently reduces the number of noise parameters to estimate ( $p$  instead of  $N$ ), and therefrom the computational cost. Secondly choosing an AR model rather than a more general class such as autoregressive-moving average models (ARMA) allows us to easily estimate the parameters from the discontinuous data, while ARMA models usually involve computationally expensive optimization procedures, or direct estimation of the autocovariance function which is not accurate when data are missing. Furthermore any moving-average model can be approximated by a high order AR model as discussed by Durbin [22].

The AR parameters  $\boldsymbol{\theta} = (a_i, \sigma^2)$  are estimated thanks to Burg's algorithm adapted to the missing data case [17]. This technique relies on the minimization of forward and backward residuals of the model (14) through a recursive procedure that increases the order  $k$  of the AR model at each step, until  $k$  reaches  $p$ . This algorithm takes advantage of all segments of available data. For a given order  $k$ , only the segments of size  $N_s > k$  can be used for the estimation. Note that the proper AR order must be previously determined according to some criteria such as Akaike's, as discussed later.

For the first iteration, the AR estimation is performed on the residuals of the OLS estimation  $\hat{\mathbf{z}}_o = \mathbf{y}_o - A_o \hat{\boldsymbol{\beta}}_{\text{OLS}}$  instead of  $\mathbf{y}_o$ , where  $\hat{\boldsymbol{\beta}}_{\text{OLS}}$  is the result of the simple estimate given by Eq. (8). This reduces the disturbance of deterministic components onto the estimation of the noise parameters.

### B. Step 2: computation of the weighted vectors with the Kalman filter

The determination of the AR parameters gives access to the noise autocovariance function. The aim of this step is to use this result to calculate the weighted observation vector  $L_o^{-1}\mathbf{y}_o$  and weighted model matrix  $L_o^{-1}A_o$  involved in the expression of the estimator (13). The matrix  $L_o$  indirectly depends on the AR parameters via the autocovariance function, since:

$$\Sigma_o(m, l) = R_z(|n_m - n_l|) \quad \forall (m, l) \in \llbracket 0, N_o - 1 \rrbracket^2, \quad (16)$$

where the autocovariance function  $R_z$  is estimated by taking the inverse Fourier transform of Eq. (15).

Unlike the case of complete stationary random series, the observed data in a missing pattern do not have a circulant nor Toeplitz correlation matrix, because the  $n_i$ 's are not regularly arranged. Therefore the matrix  $\Sigma_o$  cannot be inverted by efficient techniques such as Levinson or FFT algorithms. If the data sample is large (like in the MICROSCOPE case where typically  $N \sim 10^6$ ), this creates memory difficulties to store such a matrix. That is why we present a way of avoiding the direct inversion using a Kalman algorithm to compute the weighted data.

The relationship between GLS and Kalman filtering is explained as follows. Following the notation of Gómez and Maravall [20], an AR process can be described by the state-space representation:

$$\mathbf{x}(n) = F\mathbf{x}(n-1) + G\epsilon(n), \quad (17)$$

$$z_n = H^T \mathbf{x}(n). \quad (18)$$

The above equations are the state equation and the observation equation of the Kalman Filter.  $\mathbf{x}(n)$  is the state vector at time  $n$ , defined by:

$$\mathbf{x}(n) \equiv (z_n \ z_{n+1|n} \ \dots \ z_{n+p-1|n})^T,$$

where  $z_{n+k|n}$  is the conditional expectation of  $z_{n+k}$  given the observations before time  $n$ .  $H$  is the matrix linking the state vector to the observations, and simply reads  $H = (1 \ 0 \ \dots \ 0)^T$ . The model matrix  $F$  and the model noise vector  $G$  are calculated from the AR parameters and are defined in appendix C.

The Kalman filter aims at calculating the *a priori* estimate of the state vector along with its variance at each time  $n$  given all the observations until time  $n-1$ , that is:

$$z_{n|n-1} \equiv \mathbb{E}[z_n | z_0, z_1, \dots, z_{n-1}], \quad (19)$$

$$\sigma_{n|n-1}^2 \equiv \text{Var}[z_n | z_0, z_1, \dots, z_{n-1}]. \quad (20)$$

We define the normalized innovation vector  $\mathbf{e}$  whose elements are calculated with the Kalman residuals and their standard errors:

$$e_n \equiv (z_n - z_{n|n-1}) / \sigma_{n|n-1}. \quad (21)$$

Since  $z_{n|n-1}$  is actually the projection of  $z_n$  onto the subspace generated by  $(z_0 \ \dots \ z_{n-1})$ , Eq. (21) is equivalent

to a Gram-Schmidt orthogonalization procedure. As a result, the  $e_n$ 's are uncorrelated. In addition,  $z_{n|n-1}$  is a linear combination of  $z_i$ ,  $i < n$ , thus the normalized innovation vector can be expressed as:

$$\mathbf{e} = T\mathbf{z}, \quad (22)$$

where  $T$  is a lower triangular matrix with diagonal elements equal to one. If we calculate the autocovariance of Eq. (22) we find that

$$\text{Cov}[\mathbf{e}] = T\Sigma_z T^\dagger \Rightarrow \Sigma_z = (TT^\dagger)^{-1}, \quad (23)$$

where the implication is based on the fact that  $\text{Cov}[\mathbf{e}]$  is equal to the identity matrix. This last equation shows that the matrix  $T$  is equal or proportional to the inverse of the Cholesky decomposition  $L^{-1}$  of the covariance matrix  $\Sigma_z$ . However, in our problem this is not exactly true. The derived equalities are only valid if the random data truly follows the AR process, which is not the case in our approach since the AR model is just an approximation of the real underlying random process. We thus assume that the Kalman output  $\mathbf{e}$  is only approximately equal to  $L^{-1}\mathbf{z}$ .

If data are missing, the classic Kalman procedure must be slightly modified to properly deal with missing data, as explained by Jones [18], but the components of the normalized innovation vector  $\mathbf{e}$  corresponding to missing data are ignored in the estimation at step 3.

### C. Step 3: computation of the GLS estimate

In the previous paragraph we showed how to perform a quasi orthogonalization of the observation vector, which is exactly what is needed to compute an approximate version of the Generalized Least Squares (GLS) estimate.

The estimator in Eq. (13) can be rewritten:

$$\hat{\boldsymbol{\beta}} = (E_o^\dagger E_o)^{-1} \cdot E_o^\dagger \mathbf{e}_o, \quad (24)$$

where, with obvious notation, we denote the normalized innovation vectors  $\mathbf{e}_o = T_o \mathbf{y}_o$  and  $E_o = T_o A_o$ , calculated with the outputs of the Kalman filter algorithm, respectively applied to the observed signal and to each columns of the model matrix. Both vectors are obtained by keeping elements corresponding to observed data only. The Kalman algorithm is thus used here as a device to compute the weighted vectors involved in the GLS.

## IV. THEORETICAL UNCERTAINTY AND DETECTION ISSUES

This is of key interest to be able to assess the statistical uncertainty of a given estimation, especially in a context where the experiment cannot be reproduced a large number of times. In this section we present a tool to quantify the uncertainty of the regression result and to

give a confidence threshold for the detection of the signal of interest. To achieve this goal, the estimator variance matrix must be estimated.

The correlation matrix can be approximated under the assumption that the AR model is a good approximation of the real noise correlations. This hypothesis is equivalent to assuming that the estimator has minimal variance (*i.e.* that the estimator is the BLUE). Let  $C$  be the covariance matrix of the estimator  $\hat{\boldsymbol{\beta}}$ . Then Eq. (9) gives, by replacing  $W$  by  $T_o$  and  $A$  by  $A_o$ :

$$\hat{C} \approx \sigma_0^2 \left( E_o^\dagger E_o \right)^{-1}, \quad (25)$$

where  $\sigma_0$  accounts for the fact that the covariance is known up to a proportionality constant. For an unbiased estimator (*i.e.* the model matrix  $A_o$  describes all the deterministic components of the signal) this can be estimated by:

$$\hat{\sigma}_0^2 = \frac{\hat{\mathbf{e}}_z^\dagger \hat{\mathbf{e}}_z}{N_o - q}, \quad (26)$$

where  $\hat{\mathbf{e}}_z$  is the vector of weighted residuals defined by  $\hat{\mathbf{e}}_z \equiv \mathbf{e}_o - E_o \hat{\boldsymbol{\beta}}$ . The statistic to be considered is:

$$Z_k \equiv \frac{\hat{\beta}_k}{\sqrt{\hat{C}_{k,k}}}, \quad (27)$$

where  $k$  is the index corresponding to the parameter of interest in the vector  $\boldsymbol{\beta}$ . For our application  $\beta_k$  is the EP parameter  $\delta$ . Here we assume that the underlying process is Gaussian, which is reasonable in the case of the MICROSCOPE instrumental noise. Then under the assumption that there is no violation signal (hypothesis  $H_0$ ),  $Z$  approximately follows a Normal law with mean zero and unit variance. A detection threshold with a  $(1 - \alpha)$ -confidence level is given by imposing that the probability to observe a value above the threshold, under  $H_0$ , must be lower than  $\alpha$ . This gives the threshold  $z = \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right)$ , where  $\Phi$  is the Normal Cumulative Distribution Function (CDF). Therefore if  $|Z|$  is above the threshold, then a signal is detected with a confidence of  $100(1 - \alpha)\%$ . Conversely, for a given estimation of the EP signal, the violation can be claimed at a  $100(2\Phi(Z) - 1)\%$  confidence level. Typically, a 99% confidence test requires  $z = 2.56$ .

## V. APPLICATION TO SIMULATED DATA OF THE MICROSCOPE MISSION

### A. The MICROSCOPE experiment

The Weak Equivalence Principle (WEP) is at the basis of General Relativity. Its concrete manifestation is the Universality of Free Fall, stating that a body in a gravitational potential is accelerated independently of its mass

and internal composition. Current efforts to build new unification theories may call this principle into question [23], postulating the existence of additional fundamental interactions. To provide an experimental discrimination of these theories, the goal of the MICROSCOPE space mission is to test the WEP within a precision of about  $10^{-15}$  never reached by previous ground-based experiments [24, 25]. This space-borne test takes advantage of the duration of the fall by integrating the data over several orbits.

The mission payload is an ensemble of two electrostatic differential accelerometers composed by a cage containing two cylindrical and co-axial test-masses (TM). One accelerometer is devoted to the EP test, while the other serves as a reference. In the first accelerometer, the two TM have different compositions: one is made of Platinum Rhodium alloy (PtRh) and the other of Titanium alloy (TA6V) [26]. In the second accelerometer the TM are both made with PtRh. The masses, whose potential is kept constant via a thin gold wire, are servo-controlled by a set of electrodes to follow the same trajectory. The MICROSCOPE science signal is the difference between the accelerations applied to the two TM, which are deduced from the applied electrostatic forces needed to maintain them relatively motionless at the center of the cage. A drag-free system ensures that the measured common acceleration, *i.e.* the mean acceleration of the two TM, is nullified. A violation of the WEP would result in a difference between the two measured accelerations.

The violation signal is expected to be periodic with a frequency  $f_{EP}$  because of the projection of the gravitational acceleration onto the science axis of the instrument during the orbital trajectory. For a satellite inertial pointing session,  $f_{EP}$  is equal to the orbital frequency. For a slowly rotating satellite in the orbital plane, this is equal to the sum of the orbital frequency and the satellite spin frequency. The duration of each session is chosen in order to reach a standard deviation error of about  $10^{-15}$  on the EP parameter  $\delta$ , which is almost equal to the Eötvös parameter. The inertial and spin sessions last respectively 120 and 20 orbits. The specificity of the data samples to be analyzed in the MICROSCOPE mission is that the signal of interest has a low signal-to-noise ratio (SNR) that lies at low frequencies ( $10^{-4}$  -  $10^{-3}$  Hz) in a time series with a broad frequency range ( $10^{-5}$  - 2 Hz), blurred by a colored noise containing most of its power in higher frequencies (above  $10^{-1}$  Hz). In addition, long time series must be analyzed to achieve a sufficient SNR, including about  $5 \times 10^5$  data points for a spin session.

## B. Considered data sets

We apply the KARMA method to a time series simulated with a mission simulator. The simulation output is the differential acceleration vector  $\gamma$  equal to the acceleration difference between the two masses. This time series is sampled at  $f_s = 4$  Hz and lasts 20 orbits. This

corresponds to a spin session, for which the orbital frequency is equal to  $1.7 \times 10^{-4}$  Hz and the spin frequency is  $7.7 \times 10^{-4}$  Hz. The EP frequency is then equal to the sum  $f_{EP} = 9.4 \times 10^{-4}$  Hz.

In addition to the signal of interest, other perturbations are present in the measurement as indicated in Eq. (1). They are mainly due to gradient terms between the center of mass of the two TM, the relative motion of the TM, and coupling with the common mode because of instrument defects. During the experiment, the instrument or the satellite undergoes excitations that favor the SNR to measure their amplitudes  $\alpha_i$ . The corresponding accelerations  $s_{p,i}$  are either modeled or measured, such that the perturbations can be subtracted from Eq. (1).

Nevertheless, in this simulation we allow for the presence of gradient perturbations. They come from the slight off-centering of the test-masses, leading to gravity and inertia gradient terms. In the simulation we assume that the TM are off-centered by 20 microns along the  $x$  and  $z$ -axis which are in the orbital plane. Note that although an off-centering along the  $y$ -axis can also exist, it is estimated by means of dedicated calibration sessions and corrected numerically before the EP estimation. The EP parameter is simulated at a level of  $3 \times 10^{-15}$ . Thus the regression model  $A$  contains the true acceleration signal  $g_x(t)$ , to which we add the two perturbations modeled with our knowledge of gravity and inertia gradients. The noise added to the data is generated from the PSD model given by Eq. (6). The signal model reads:

$$\gamma(t) = \delta g_x(t) + \Delta_x T_{xx}(t) + \Delta_z T_{xz}(t) + z(t), \quad (28)$$

where we have noted  $g_x$  the gravitational acceleration projected onto the  $x$ -axis,  $T_{ij}$  the components of the gradient tensor, and  $\Delta_i$  the off-centerings. Thus in this case there are three regression parameters:  $\delta$ ,  $\Delta_x$  and  $\Delta_z$ .

We consider two types of gap pattern. The first one is a “tank crackle type” window  $w_a$  that is generated so that all gaps are of equal duration (0.5 second) and their positions are randomly distributed on the sample (uniform distribution with 260 gaps per orbit). The second one is a “telemetry losses type” window  $w_b$  where the gaps durations are drawn from a distribution similar to the telemetry thread of the PICARD mission [27], with a standard duration of one minute. Their positions are distributed in the same way as for the first window. Each window represents the same fraction of missing data, of about 2%. Thus window  $w_a$  comprises more gaps than window  $w_b$  (larger  $N_g$ ) but gaps are shorter in average. To illustrate this, we plot in Fig. 3 an extract of the time series where the data interruptions of each window are identified by vertical grey bars.

We apply the KARMA method and compare the result to the Ordinary Least Squares estimate with missing data to assess the improvement. We also compare the result to the reference given by the OLS estimator in the case without gaps.



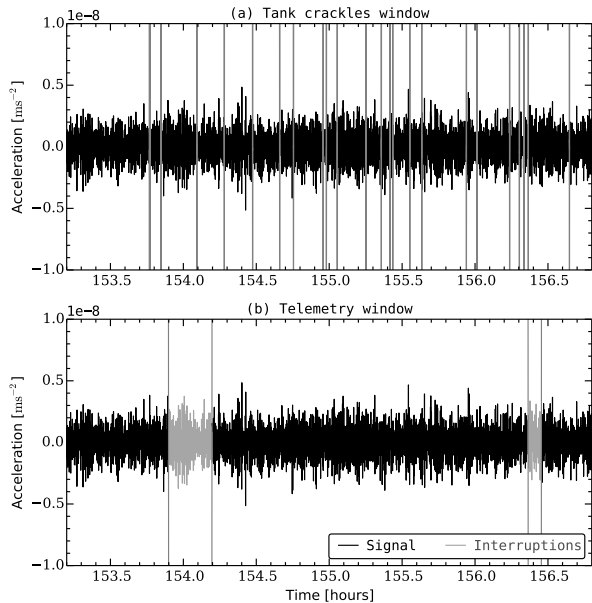


FIG. 3. Fraction of the temporal series (black) with the interruption times represented by the grey vertical lines for the two windows  $w_a$  (top) and  $w_b$  (bottom).

### C. PSD estimate from the AR fit

Before starting the whole process, the order of the AR model must be chosen at step 1. The choice of the order depends on the PSD of the noise affecting the measurement, and on the observation window. A way to properly choose the order  $p$  is to minimize the Akaike's Information Criterion [28] defined as  $AIC(p) = 2p - 2 \log(L_{\max}(p))$ , where  $L_{\max}$  is the maximized log-likelihood. In the case of an AR model, this can be expressed in terms of the estimate of the AR residual variance  $\hat{\sigma}^2$  which is directly computable from the residuals of the Burg's algorithm:

$$AIC(p) = 2p + N_o \log(\hat{\sigma}^2). \quad (29)$$

Applying Burg's algorithm to the residual series  $\hat{z}$  defined in section III A with increasing order  $k$  allows us to find the order that minimizes the AIC.

In the MICROSCOPE case,  $AIC(p)$  is an asymptotically monotonic decreasing function. In this configuration one possibility is to choose the order  $p$  from which there is no significant improvement in the AIC, *i.e.* the minimum order where the AIC is close enough to the asymptote. For the studied noise, the AIC typically reaches a plateau from  $p = 200$ .

Nevertheless, in case of very frequent missing data (e.g. tank crackles), the variance of the AR coefficients estimates increases with the order, and so does the variance of the AIC. This is due to the decrease of the number of usable data segments (with length larger than  $p$ ). This

can lead to overestimating the optimal order  $p$ . To overcome this difficulty we can modify the AIC criterion as suggested by Bos *et al.* [29] by introducing a penalty accounting for the increasing estimation variance. We choose to replace  $N_o$  by  $p \left( \sum_{i=1}^p \frac{1}{N_i} \right)^{-1}$  where  $N_i$  is the number of usable segments to estimate the coefficient  $a_i$ . When applying this criterion to our simulation with window  $w_a$ , we find an optimal order of  $p = 60$ .

The process converges after 2 iterations, because the first estimate of the PSD is influenced by the high amplitude perturbations of the gradient terms: the main peak has an amplitude of  $2.4 \times 10^{-11} \text{ ms}^{-2}$  at  $2f_{EP}$ , and other peaks are present at  $f_{orb}$  and  $2f_{orb}$ . In comparison, the EP violation corresponds to an amplitude of  $1.2 \times 10^{-14} \text{ ms}^{-2}$ .

We plot in Fig. 4 the estimate of the PSD (red curve) obtained with the AR coefficients calculated by Burg's algorithm with the tank crackles (a) and telemetry (b) windows, along with the real PSD (black curve). The level of noise of the masked data is shown by the black dotted line. In addition to the selected order  $p = 60$ , we also show the AR spectrum estimate made with a larger order ( $p = 200$  with window  $w_a$ , and  $p = 2000$  with window  $w_b$ ) to illustrate the effect of  $p$ . In both cases, we see that the overall shape of the PSD is well described by the AR model, especially the  $f^4$  slope. However, there is a bias which increases as the frequency decreases. The reasons why the AR model cannot accurately describe the low frequency PSD are two-fold:

1. The order of the AR model is finite, and limited by the longest segment of consecutive available data (this is typically 700 for the window  $w_a$ , and 5000 for  $w_b$ ). Given that AR models cannot describe  $1/f$  spectra with a finite number of parameters, a larger order is necessary to reduce the bias (and the bias is zero when  $p$  tends to infinity).
2. In the Burg estimation procedure, the larger the AR order, the larger the variance of the AR coefficients estimates, because there are fewer segments of corresponding lengths. This is why we do not choose the highest possible order, for which segments of corresponding length are rare.

Since window  $w_b$  has more spaced and longer gaps than window  $w_a$ , it allows for a higher possible AR order leading to a better restitution of the low frequency shape of the PSD, with a reasonable variance (see Fig. 4). However we choose  $p = 60$  even in the case of window  $w_b$  for computational reasons, given that this is the high frequency restitution that matters for a parameter regression purpose, as we shall see in the next paragraph.

### D. Regression results

The results of the linear regression are summarized in Table I, with  $p = 60$ . In order to test the precision of our

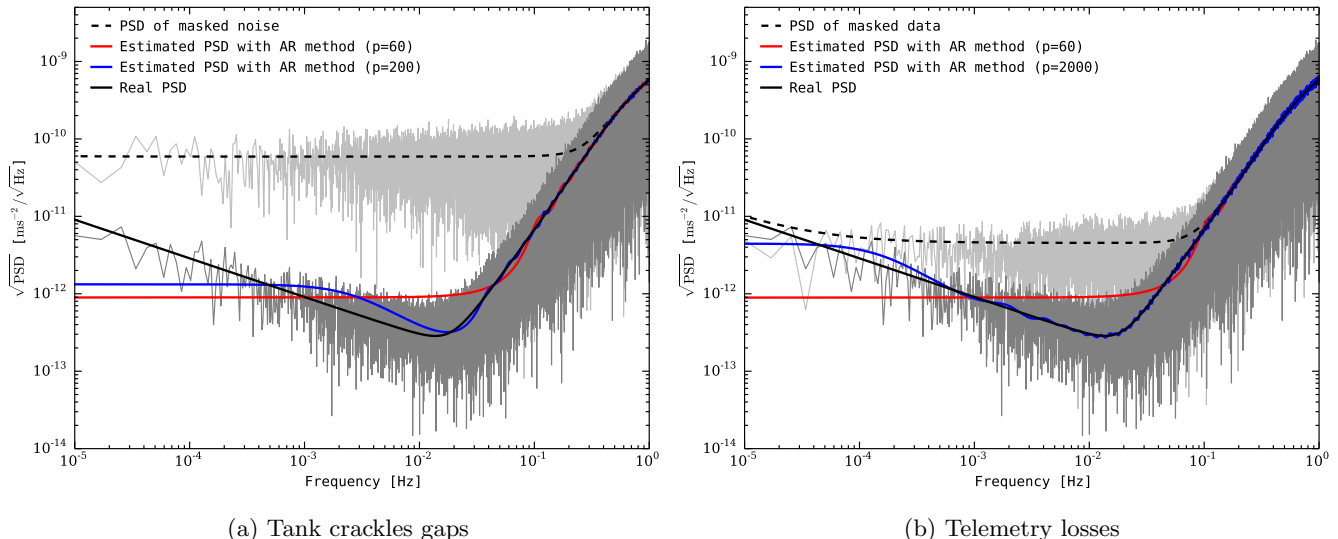


FIG. 4. PSD estimates of the noise in presence of missing data. The black dashed curve is an estimate of the masked data PSD [obtained using Eq. (5)], the black solid curve is the actual noise PSD, and the red and blue curves are the PSD estimates of the AR model obtained with Burg’s algorithm. The periodograms of the regression residuals are also plotted for the complete (dark grey) and masked (light grey) cases.

method, we have drawn 400 realizations of the noise and run our estimation algorithm for each of them, as well as the OLS estimator. The number of draws is chosen such that the error on the true value of the standard deviation of the EP parameter does not exceed  $10^{-16}$  with a 99% confidence.

The third column of the table indicates the true value of the parameters. Columns 4 to 6 show the performance of the OLS estimator: the sample average  $\hat{\mu}$ , the theoretical standard deviation  $\sigma_{th}$  given by Eq. (9) and calculated with the real PSD, and the sample standard deviation of the 400 estimates. The last three columns show the results obtained with the KARMA method, and are detailed below.

The sample mean  $\hat{\mu}$  of the estimates obtained with the KARMA method converges to the true value of the parameters (seventh column of table I), showing that the constructed estimator is unbiased.

We also calculate the sample standard deviation of the EP parameter. For short and numerous gaps (tank crackles window) we find  $\hat{\sigma} = 1.1 \times 10^{-15}$  with our method instead of  $6.5 \times 10^{-14}$  with the OLS estimator. Thus our method enables us to divide the stochastic error by a factor 60 with respect to the OLS.

For fewer and longer gaps (telemetry window), we find  $\hat{\sigma} = 9.8 \times 10^{-16}$  instead of  $5.1 \times 10^{-15}$  with the OLS. We notice that such a gap pattern has less impact on the estimation performance, because it leads to a lower frequency leakage as confirmed by Eqs. (A3) and (A4) of appendix A (also see Bergé *et al.*, in prep.).

These are satisfying results since the theoretical uncertainty of the OLS without any missing data is equal to  $9.6 \times 10^{-16}$ . The detection test is positive with a confidence greater than 99% in both cases.

The improvement is also significant for the other parameters. Even if they are already well estimated by the OLS, their uncertainty is reduced by almost two orders of magnitude for the tank crackles window.

For each draw, we estimate the uncertainty  $\hat{\sigma}_{AR}$  using the approximate formula (25). We then calculate the sample average of this estimate over the 400 draws, and record the results in the table. We find  $1.2 \times 10^{-15}$  for window  $\mathbf{w}_a$  and  $9.3 \times 10^{-16}$  for window  $\mathbf{w}_b$ . This is close to the calculated sample standard deviation, meaning that when having only one realization at hand, one can estimate the error with an acceptable accuracy. The estimated error does not vary much from one estimation to another, and stays within an interval of  $\pm 10^{-16}$  around the mean.

The estimate  $\hat{\sigma}_{AR}$  of the real regression error may be biased, depending on the frequency of the estimated signal. This can be explained by Fig. 4, where we observe that the PSD of the AR model is biased at low frequency. As a result, the lower the signal frequency, the larger the bias on the estimated variance. This is particularly true around zero, where the AR PSD is below the real one. However, the overall shape of the real PSD is well captured by the AR model, which is enough to cancel the leakage due to the window and get a precision of  $1 \times 10^{-15}$  for the EP estimation, in agreement with the mission requirement.

## VI. CONCLUSION AND DISCUSSION

We have shown that the presence of gaps in time series affected by correlated noise has a strong impact on the classical Fourier analysis and on the precision of the

TABLE I. Mean and standard deviations on the estimation of the parameters of interest using OLS and the KARMA method. In both cases we present (from left to right) the estimation average calculated on a sample of 400 estimates, the analytical standard deviation, and the sample standard deviation. For OLS, the analytical uncertainty  $\sigma_{th}$  is given by Eq. (9), which is exact. For the KARMA method,  $\hat{\sigma}_{AR}$  is the average of the uncertainties estimated for each draw with Eq. (25).

Window	Param.	True	Ordinary Least Squares			Kalman-AR Model Analysis		
			$\hat{\mu}$	$\sigma_{th}$	$\hat{\sigma}$	$\hat{\mu}$	$\hat{\sigma}_{AR}$	$\hat{\sigma}$
Complete data	$\delta [10^{-15}]$	3	3.01	0.96	1.02	2.98	0.92	0.96
	$\Delta_x [\mu\text{m}]$	20	20.0	0.003	0.005	20.0	0.004	0.003
	$\Delta_z [\mu\text{m}]$	20	20.0	0.003	0.005	20.0	0.004	0.003
Tank crackles	$\delta [10^{-15}]$	3	8.82	62.3	65.2	2.98	1.19	1.14
	$\Delta_x [\mu\text{m}]$	20	20.0	0.290	0.296	20.0	0.006	0.004
	$\Delta_z [\mu\text{m}]$	20	20.0	0.292	0.314	20.0	0.006	0.005
Telemetry	$\delta [10^{-15}]$	3	3.15	5.20	5.07	2.98	0.93	0.98
	$\Delta_x [\mu\text{m}]$	20	20.0	0.021	0.021	20.0	0.004	0.003
	$\Delta_z [\mu\text{m}]$	20	20.0	0.024	0.024	20.0	0.004	0.003

ordinary least squares fits of harmonic signals. This is due to the frequency leakage of the noise power, which can increase the uncertainty of the fit by several orders of magnitude, even if the percentage of missing data is small.

We proposed a method that we dubbed ‘‘KARMA’’, which provides a general way to perform precise linear regressions with large and incomplete data sets affected by unknown colored noise, and that we applied to mock MICROSCOPE data. The estimation variance is decreased down to the same order of magnitude as the least squares estimator with full data, altered by the natural loss of signal due to the  $1/\sqrt{N}$  dependence. The method tends to approach the minimum variance estimator of the available data, by approximating the noise autocovariance with a high order AR model.

Our method uses a weighting of the data relying on the estimation of the shape of the PSD. As a result, the performance of the regression mainly depends on the ability of the autoregressive PSD estimate to recover the part of the spectrum that is responsible for the leakage, which is the high frequency part increasing as  $f^4$  in the MICROSCOPE case. Although this is not shown here, the method has also been successfully tested in a case where the leaking power comes from a thermal  $1/f^2$  noise projected onto high frequencies. The AR PSD then accurately fits the low frequency slope and allows us to improve the possible regression of high frequency components.

In addition, the outputs allow us to evaluate the variance of the estimator from a single estimation. We recover the magnitude of the true precision, equal to  $10^{-15}$  in our MICROSCOPE illustration. The variance is not estimated with a better accuracy because of the low frequency bias of the AR PSD estimator. This bias depends on the missing data pattern, and more particularly on the length of the longest uninterrupted data segment, as well as the number of long segments. This determines the AR order to be chosen, resulting in a trade-off between the

bias and the variance of the PSD estimate.

Concerning the scientific objective of the MICROSCOPE mission, the above discussion demonstrates that based on the current noise model of the accelerometers, we will be able to get a 99% (resp. 68%)-confidence level detection of a  $3 \times 10^{-15}$  (resp.  $1 \times 10^{-15}$ ) EP violation signal, even in the presence of missing data, for a 20 orbit-measurement session (completed in 1.4 days). The mission should include more than 70 sessions of this type, allowing for a detection at the 99% level even for an amplitude of  $1 \times 10^{-15}$ . This has been done for short and very frequent gaps to represent acceleration peaks or saturations due to MLI or tank crackles, as well as for longer and fewer gaps to simulate telemetry interruptions.

Further developments will concentrate on how to increase the accuracy of the noise PSD estimate, for example by using the AR model to perform missing data imputation. Indeed, although this is computationally more expensive, the AR model can be exploited in a Gaussian process regression approach [30] to estimate the missing values.

Finally, there are two potential limitations to the presented method that can be addressed in further extensions. On the one hand, although the AR model can be a good approximation to any PSD and can be fitted very efficiently, it is still a parametric and thus restrictive model. On the other hand, noise and signal parameters are estimated iteratively but separately, so that each step is done conditionally to the previous one. This may result in a loss of accuracy. As a result, a possible generalization is to use the proposed method as an efficient initialization procedure for a more general regression algorithm that would maximize the full likelihood without any prior noise model.

### Appendix A: PSD deformation in the case of random missing data patterns

We derive here the PSD of the masked data in the case where the gaps positions in the time series are drawn from a uniform distribution. Let  $N$  be the length of the time series,  $N_g$  the number of gaps, and  $n_{b,i}$  the indices indicating the location of the beginning of each gap (such that  $w_{n_{b,i}} = 0$ ). Each gap ends at the location  $n_{b,i} + dn_i$  (we adopt the convention  $w_{n_{b,i}+dn_i} = 1$ ). By uniformly distributed, we mean that  $n_b$  is a random variable following a discrete uniform distribution on the interval  $\llbracket 0, N-1 \rrbracket$ . We also allow the gap duration  $dn$  to be randomly distributed. The window vector is then generated by drawing  $N_g$  realizations of  $n_b$  and  $dn$ .

The probability  $P$  to observe a data at a time  $n$  is calculated as follows:

$$\begin{aligned} P &= \text{P}(w_n = 1) \\ &= \prod_{i=0}^{N_g-1} \text{P}(n < n_{b,i} \text{ or } n \geq n_{b,i} + dn_i) \\ &= [1 - \text{P}(n_b \leq n) + \text{P}(n_b + dn \leq n)]^{N_g}. \end{aligned} \quad (\text{A1})$$

The cumulative probability function of  $n_b$  is given by:

$$\text{P}(n_b \leq n) = \frac{n+1}{N}. \quad (\text{A2})$$

In the case where the duration of the gaps is fixed (*i.e.*  $dn_i = dn_0 \forall i$ ), Eq. (A1) gives:

$$\begin{aligned} \text{P}(w_n = 1) &= \left[ 1 - \frac{n+1}{N} + \frac{n-dn_0+1}{N} \right]^{N_g} \\ &= \left[ \frac{N-dn_0}{N} \right]^{N_g}. \end{aligned} \quad (\text{A3})$$

Therefore the probability law of  $w_n$  is a Bernoulli's law of parameter  $P$ . Its expectation is  $\mu_w = P$  and its variance is  $\sigma_w^2 = P(1-P)$ . We notice that  $P$  is independent of time, and the autocovariance function of  $\mathbf{w}$  is simply  $R_w(n) = \sigma_w^2 \delta(n)$  where  $\delta(n)$  is the delta Dirac function. Then we use Eq. (5) to calculate the PSD of the masked data:

$$S_y(f) = P^2 \cdot S_z(f) + P(1-P) \int_{-\frac{f_s}{2}}^{\frac{f_s}{2}} S_z(f') df'. \quad (\text{A4})$$

### Appendix B: Derivation of a simplified equation for the OLS variance in the harmonic case

We derive here the approximate expression of the variance of the ordinary least squares estimator used in Section II.

We start from Eq. (9). In the case of a simple harmonic model, the matrix  $A_w$  is a column matrix and the

covariance formula can be written as:

$$\text{Var}(\hat{\delta}) = \frac{A_w^\dagger \Sigma_w A_w}{(A_w^\dagger A_w)^2}.$$

As reminded in Eq. (10), the covariance matrix is diagonalizable in the Fourier space. We keep the same notations in the following. In addition, we use the fact that the Discrete Fourier Transform (DFT) operator is a Vandermonde matrix (since  $M^\dagger M = NI$  with  $I$  the identity matrix), therefore the variance can be rewritten in terms of the DFT of the windowed model  $A_w$ , noted  $\tilde{A}_w = MA_w$ :

$$\text{Var}(\hat{\delta}) = N f_s \frac{\tilde{A}_w^\dagger D \tilde{A}_w}{(\tilde{A}_w^\dagger \tilde{A}_w)^2}.$$

By developing this expression we get:

$$\text{Var}(\hat{\delta}) = \frac{\sum_{k=0}^{N-1} |\tilde{A}_{wk}|^2 N f_s \hat{S}_{y_k}}{\left( \sum_{k=0}^{N-1} |\tilde{A}_{wk}|^2 \right)^2}.$$

For the windowed harmonic model  $A_{wn} = w_n \gamma_{\text{EP}} \sin(2\pi n f_{\text{EP}}/f_s + \phi_{\text{EP}})$ ,  $\tilde{A}_w$  is the convolution of the DFT of the window and the DFT of the EP signal. In the case of a random window,  $|\tilde{A}_w|$  usually peaks at the EP frequency with a value of  $\gamma_{\text{EP}} N_o/2$  where  $N_o$  is the number of observed data (where  $w_n = 1$ ). To simplify the calculations, we neglect the terms at all other frequencies. This amounts to ignoring the leakage of the harmonic signal (but note that the leakage of the noise component is present in  $S_y$  through equation 5). Furthermore, if we assume that the integration period is an integer multiple of the EP period (*i.e.* there exist an integer  $k_{\text{EP}}$  such that  $f_{\text{EP}} = k_{\text{EP}} f_s/N$ ), then we have:

$$\text{Var}(\hat{\delta}) \approx \frac{\gamma_{\text{EP}}^2 \frac{N_o^2}{4} N f_s (S_y(f_{\text{EP}}) + S_y(-f_{\text{EP}}))}{\left( 2 \times \gamma_{\text{EP}}^2 \frac{N_o^2}{4} \right)^2},$$

where we have made the approximation, valid for large  $N$ , that the DFT of the autocovariance function in Eq. (11) is equal to the real PSD. By simplifying we get equation 12:

$$\text{Var}(\hat{\delta}) \approx \frac{2 f_s N S_y(f_{\text{EP}})}{N_o^2 \gamma_{\text{EP}}^2}.$$

Note that in the case of a complete data set ( $w_n = 1 \forall n$ ) we have  $N_o = N$  and this formula is more accurate because the model  $|\tilde{A}_w|$  exactly peaks at  $\gamma_{\text{EP}} N/2$ .

### Appendix C: State space equation of an AR model

We detail here the Kalman equations presented in Section III B.

The observation matrix, the model matrix and the model noise matrix are defined respectively by:

$$\begin{aligned}
 H &\equiv (1, 0, \dots, 0)^T \\
 F &\equiv \begin{pmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 \\ -a_p & -a_{p-1} & -a_{p-2} & \dots & -a_1 \end{pmatrix} \\
 G &\equiv (1, g_1, \dots, g_{p-1})^T.
 \end{aligned}$$

The vector  $G$  whose elements are defined by  $\epsilon_n g_j \equiv z_{n+j-1|n} - z_{n+j-1|n-1}$  can be calculated from the AR parameters (see Jones [18]).

The Kalman filter equations in presence of missing data are briefly reviewed here:

### Prediction equation

$$\begin{aligned}
 \mathbf{x}(n|n-1) &= F\mathbf{x}(n-1|n-1), \\
 \Sigma(n|n-1) &= F\Sigma(n-1|n-1)F^T + Q,
 \end{aligned}$$

where  $Q \equiv GG^T$ .

### Update equation

The update equation adapted to the missing data case

can be formulated as follows:

$$\begin{aligned}
 \Sigma(n|n) &= w_n \{ \Sigma(n|n-1) - K(n)H^T \Sigma(n|n-1) \} \\
 &\quad + (1 - w_n) \{ \Sigma(n|n-1) \}, \\
 \mathbf{x}(n|n) &= w_n \{ \mathbf{x}(n|n-1) + K(n) (z(n) - H^T \mathbf{x}(n|n-1)) \} \\
 &\quad + (1 - w_n) \{ \mathbf{x}(n|n-1) \},
 \end{aligned}$$

where we defined

$$K(n) \equiv \Sigma(n|n-1)H (H^T \Sigma(n|n-1)H)^{-1}.$$

Note that if the data is not observed at time  $n$ , the state variance and the state vector are not updated and set equal to the predicted values at previous time.

## ACKNOWLEDGMENTS

The authors would like to thank all the members of the MICROSCOPE Performance team, as well as Sandrine Pires and Jean-Philippe Ovarlez, for fruitful discussions. This activity has been funded by ONERA and CNES. We also acknowledge the financial contribution of the UnivEarthS Labex program at Sorbonne Paris Cité (ANR-10-LABX-0023 and ANR-11-IDEX-0005-02).

- 
- [1] P. Touboul, G. Métris, V. Lebat, and A. Robert, *Classical and Quantum Gravity*, **29**, 184010 (2012).
- [2] J. Bergé, P. Touboul, and M. Rodrigues, *ArXiv e-prints* (2015), arXiv:1501.01644.
- [3] S. Vitale, G. Congedo, R. Dolesi, V. Ferroni, M. Hueller, D. Vetruigno, W. J. Weber, H. Audley, K. Danzmann, I. Diepholz, M. Hewitson, N. Korsakova, L. Ferraioli, F. Gibert, N. Karnesis, M. Nofrarias, H. Inchauspe, E. Plagnol, O. Jennrich, P. W. McNamara, M. Armano, J. I. Thorpe, and P. Wass, *Phys. Rev. D*, **90**, 042003 (2014).
- [4] E. Hardy, A. Levy, G. Métris, M. Rodrigues, and P. Touboul, *Space Science Reviews*, **180**, 177 (2013).
- [5] N. R. Lomb, *Astrophys. and Space Science*, **39**, 447 (1976).
- [6] J. D. Scargle, *Astrophys. J.*, **263**, 835 (1982).
- [7] D. H. Roberts, J. Lehar, and J. W. Dreher, *Astron. J.*, **93**, 968 (1987).
- [8] M. Zechmeister and M. Kürster, *Astron. Astrophys.*, **496**, 577 (2009).
- [9] S. Pires, S. Mathur, R. A. Garcia, J. Ballot, D. Stello, and K. Sato, *ArXiv e-prints* (2014), arXiv:1410.6088 [astro-ph.SR].
- [10] Y. Wang, P. Stoica, J. Li, and T. L. Marzetta, *Digital Signal Processing*, **15**, 191 (2005).
- [11] P. Stoica, J. Li, and J. Ling, *Signal Processing Letters, IEEE*, **16**, 241 (2009).
- [12] C. Röver, R. Meyer, and N. Christensen, *Classical and Quantum Gravity*, **28**, 015010 (2011).
- [13] T. B. Littenberg and N. J. Cornish, *Phys. Rev. D*, **80**, 063007 (2009).
- [14] T. B. Littenberg and N. J. Cornish, *ArXiv e-prints* (2014), 0902.0368.
- [15] D. L. Donoho and X. Huo, *Information Theory, IEEE Transactions on*, **47**, 2845 (2001).
- [16] M. Elad, J.-L. Starck, P. Querre, and D. Donoho, *Applied and Computational Harmonic Analysis*, **19**, 340 (2005).
- [17] S. De Waele and P. M. Broersen, *Signal Processing, IEEE Transactions on*, **48**, 2876 (2000).
- [18] R. H. Jones, *Technometrics*, **22**, 389 (1980).
- [19] R. Kohn and C. F. Ansley, *Biometrika*, **72**, 694 (1985).
- [20] V. Gómez and A. Maravall, *Journal of the American Statistical Association*, **89**, 611 (1994).
- [21] B. Priestley, *Spectral analysis and time series*, *Probability and mathematical statistics No. vol. 1 à 2* (Academic Press, 1982).
- [22] J. Durbin, *Biometrika*, **46**, 306 (1959).
- [23] T. Damour, F. Piazza, and G. Veneziano, *Phys.Rev.*, **D66**, 046007 (2002).
- [24] T. A. Wagner, S. Schlamminger, J. Gundlach, and E. Adelberger, *arXiv preprint arXiv:1207.2442* (2012).
- [25] C. M. Will, *Living Reviews in Relativity*, **17** (2014).
- [26] É. Hardy, A. Levy, M. Rodrigues, P. Toubou, and G. Métris, *Advances in Space Research*, **52**, 1634 (2013).
- [27] M. Rouzé, A. Hauchecorne, J. F. Hochedez, A. Irbah, M. Meftah, T. Corbard, S. Turck-Chièze, P. Boumier, S. Dewitte, and W. Schmutz, in *SpaceOps 2014* (American Institute of Aeronautics and Astronautics, 2014).
- [28] H. Akaike, *Automatic Control, IEEE Transactions on*, **19**, 716 (1974).

- [29] R. Bos, S. de Waele, and P. M. Broersen, Instrumentation and Measurement, IEEE Transactions on, **51**, 1289 (2002).
- [30] C. E. Rasmussen and C. K. I. Williams, *Gaussian processes for machine learning* (The MIT Press, 2006).