# Comparing nested data sets and objectively determining financial bubbles' inceptions

G. Demos [a],[*], D. Sornette [a],[b]

[a] *ETH Zürich, Department of Management, Technology and Economics, Zürich, Switzerland*
[b] *Swiss Finance Institute, c/o University of Geneva, Geneva, Switzerland*

## ABSTRACT

Motivated by the question of identifying the start time $\tau$ of financial bubbles, we propose an improved calibration approach for time series in which the inception of the latest regime of interest is unknown. By taking into account the tendency of a given model to overfit data, we introduce the Lagrange regularisation of the normalised sum of the squared residuals, $\chi_{np}^2(\Phi)$, to endogenously detect the optimal fitting window size $:= w^* \in [\tau : \bar{t}_2]$ that should be used for calibration, assuming a fixed pseudo present time $\bar{t}_2$. The Lagrange regularisation of $\chi_{np}^2(\Phi)$ defines the Lagrange regularised sum of the squared residuals, $\chi_\lambda^2(\Phi)$. Its performance is exemplified on a simple Linear Regression problem with a change point and compared against the performances of the Residual Sum of Squares (RSS) $:= \chi^2(\Phi)$ and RSS/(N-p) $:= \chi_{np}^2(\Phi)$, where $N$ is the sample size, $p$ is the number of degrees of freedom and $\Phi$ is the parameter vector. Applied to synthetic models of financial bubbles with a well-defined transition regime and to a number of financial time series (US S&P500, Brazil IBovespa and China SSEC Indices), $\chi_\lambda^2(\Phi)$ is found to provide well-defined reasonable determinations of the starting times for major bubbles such as the bubbles ending with the 1987 Black-Monday, the 2008 Sub-prime crisis and minor speculative bubbles on other Indexes, without any further exogenous information. The application of the method thus allows one to endogenise the determination of the starting time of bubbles, a problem that has yet not received a systematic objective solution. Moreover, the technique appears as a practical solution for comparing goodness-of-fit across unbalanced sample sizes.

© 2019 Published by Elsevier B.V.

## 1. Introduction

There is an inverse relationship between the tendency of a model to overfit data and the sample size under consideration. In other words, the smaller the sample size, the larger the flexibility for the model with a fixed number of parameters to overfit [1]. Due this characteristic feature, one cannot directly compare goodness-of-fit metrics of statistical models arbitrarily parametrised by the vector $\Phi$ of parameters such as the Residual Sum of Squares, $RSS := \chi^2(\Phi)$ or its normalised version $\frac{RSS}{(N-p)} := \chi_{np}^2(\Phi)$, over unequal sized samples. Here, $N$ denotes the sample size while $p$ is the number of degrees of freedom of a model. This has particular relevance when one is interested in selecting the optimal sub-sample of a data set to calibrate a model, a recurrent issue when estimating time series models in a moving window or/and when the model is only valid in a specific time window, which is unknown *a priori*.

Our motivation stems from the question of determining the beginning of a financial bubble conditional on a fixed end point, $\bar{t}_2$, but the technique is more generally applicable to time series exhibiting regime shifts that one is interested in detecting via the minimisation of a cost function or for comparing goodness-of-fit across unbalanced sample sizes.

There are solutions for proper model selection, such as Lasso [2] and Ridge regressions [3], where the cost function penalises large values of estimated parameters. Well-known metrics such as AIC and BIC are also standard tools for quantifying goodness-of-fit of different models [4] and for solving the trade-off between goodness-of-fit and complexity. However, results stemming from these methodologies are only comparable within the same data set and are therefore not informative when unequal sample sizes are considered.

The Statistical Process Control (SPC) and Change-Point Detection literature is very rich and suggest a number of methods for detecting "break points" on a given time series (see for example [5,6] for good reviews on the different existing methodologies). Most of the available techniques focus on quantifying the deviations of the statistical properties and of the time series features before and after the candidate change point, and qualify a change point when some score for the difference between before and after it is sufficiently large. In contrast, our method focuses on problems in which a specified process of interest is assumed to originate from some starting point, replacing a previous unknown process existing before that start time. The process of interest is a priori completely different in its statistical properties, nonlinearities and structure from the pre-existing process before the start time. The start point of the new process is the focus of our method. We thus depart from the change point detection methods based on standard assumptions on a change of some characteristics of general processes, which target a change of some moments (mean, variance, and so on) or some correlations, or graph structure and so on, between before and after the change point. Our approach, which falls in the "online" class, is based on dynamically adjusting the degeneracy of $\chi^2_{np}(\Phi)$ as a function of the number $N$ of data points in the time window ending at the "present" time. Our method seems to fill a gap in the literature regarding the proper procedure for comparing goodness-of-fit of the same model calibrated over different batches of a given data set.

Here, we propose a novel metric for calibrating endogenised start points $\bar{t}_1$, and compare nested data sets. The method empirically computes the tendency of a model to overfit a data set via what we term the "*Lagrange regulariser term*", $\lambda$. Once $\lambda$ has been estimated empirically, the cost function can be corrected accordingly as a function of sample size, giving the Lagrange regularisation of the normalised Residual Sum of Squares, $\chi^2_{\lambda}(\Phi)$. Since the resulting cost function is now convex, its minimum gives us the optimal sample size that one should use in order to calibrate a model.

We empirically test the performance of the Lagrange regularisation of the regularised Residual Sum of Squares ($\chi^2_{np}(\Phi)$), defined as $\chi^2_{\lambda}(\Phi)$, in comparison with the naive $\chi^2(\Phi)$ and $\chi^2_{np}(\Phi)$, using both linear and non-linear models as well as synthetic and real-world time-series. Our Monte Carlo simulations suggest that using the window size selected by our procedure can improve the selection of the optimal sub-sample *vis-a-vis* an *ad-hoc* choice of the window size.

Using a simple Linear Regression Monte Carlo simulation study with a time series possessing different slopes in different time segments and decorated by white noise, we show how the technique is capable of accurately detecting the transition points between segments. Moreover, motivated by the objective of detecting the beginning of financial bubbles, we also employ our novel Lagrange technique on the difficult problem [7] of calibrating the non-linear Log-Periodic Power-Law Singularity model [8–12]. Results show that we can accurately infer dates where a price exuberance regime has started [13]. Besides being a straight forward methodology for detecting optimal window sizes and comparing nested data sets – thus being an extremely convenient technique for practitioners – the technique correctly turns ill-conditioned optimisation problems into convex ones for different levels of noise and across a number of different time series.

This paper is structured as follows. Section 2 explains the motivation behind the proposed Lagrange regularising term. Moreover, we provide details of the derivation of the Lagrange regulariser term $\lambda$ as well as the analytical expression for computing the tendency of a model to overfit data. In Section 3, we make use of a simple OLS regression to test the empirical performance of the Lagrange regularisation of $\chi^2_{np}(\Phi)$ on the problem of optimal sub-sample selection. Section 4 shows how the regulariser can be used alongside with the LPPLS model of financial bubbles [8] in order to diagnose the beginning of financial bubbles. Empirical findings are given in Sections 5 and 4.2 concludes.

## 2. Formulation of calibration with varying window sizes: How to endogenise $t_1$ and make different window sizes comparable

In this section, we present the motivation behind our novel technology as well as its building blocks.

Consider the time window $[t_1+1, t_2]$ containing $t_2-t_1$ time points on which a measurable process is defined. Calibrating the realisation of this process in the time window $[t_1, t_2]$, one obtains the normalised mean-squared residuals, defined as the sum of squares of the residuals divided by the number $t_2 - t_1$ of points in the sum corrected by the number of degrees of freedom $p$ of the model,

$$\chi^2_{np}(\Phi) := \frac{1}{(t_2 - t_1) - p} \sum_{i=t_1}^{t_2} r_i(\Phi)^2 \, , \tag{1}$$

with

$$r_i(\Phi) = y_i^{data} - y_i^{model}(\Phi) \, , \tag{2}$$

where $\boldsymbol{\Phi}$ denotes the set of model parameters to fit including *a priori* the left end point $t_1$ of the calibration window. The term $y_i^{model}(\boldsymbol{\Phi})$ corresponds to the theoretical model and $y_i^{data}$ is the empirical value of the time-series at time $i$.

For a fixed right end point $t_2$ of the calibration window, we are interested in comparing the results of the fit of the model to the empirical data for various left end points $t_1$ of the calibration window. The standard approach assumes a fixed calibration window $[t_1 + 1, t_2]$ with $N = t_2 - t_1$ data points. In order to relate the two problems, we consider the minimisation of $\chi_{np}^2(\boldsymbol{\Phi})$ at fixed $t_2 - t_1$ (for a fixed $t_2$) as minimising a general problem involving $t_1$ as a fit parameter augmented by the condition that $t_2 - t_1 = N$ is fixed. This reads

$$\text{Min } \chi_\lambda^2(\boldsymbol{\Phi}) , \tag{3}$$

with

$$\chi_\lambda^2(\boldsymbol{\Phi}) := \frac{1}{(t_2 - t_1) - p} \sum_{i=t_1}^{t_2} r_i(\boldsymbol{\Phi})^2 + \lambda(t_2 - t_1) , \tag{4}$$

where we have introduced the Lagrange parameter $\lambda$, which is conjugate to the constraint $t_2 - t_1 = N$. Once the parameters $\boldsymbol{\Phi}$ are determined, $\lambda$ is obtained by the condition that the constraint $t_2 - t_1 = N$ is verified.

Since data points are discrete, the minimisation of (4) with respect to $t_1$ reads

$$
\begin{aligned}
0 = \quad & \chi_\lambda^2(\boldsymbol{\Phi})(t_1 + 1) - \chi_\lambda^2(\boldsymbol{\Phi})(t_1) = \frac{1}{(t_2 - t_1 - p - 1)} \sum_{i=t_1+1}^{t_2} r_i(\boldsymbol{\Phi})^2 - \frac{1}{t_2 - t_1 - p} \sum_{i=t_1}^{t_2} r_i(\boldsymbol{\Phi})^2 - \lambda \\
= \quad & \frac{1}{t_2 - t_1 - p} \left( 1 + \frac{1}{t_2 - t_1 - p} + \mathcal{O}\left( \frac{1}{(t_2 - t_1 - p)^2} \right) \right) \sum_{i=t_1+1}^{t_2} r_i(\boldsymbol{\Phi})^2 - \frac{1}{t_2 - t_1 - p} \sum_{i=t_1}^{t_2} r_i(\boldsymbol{\Phi})^2 - \lambda , \\
= \quad & -\frac{1}{t_2 - t_1 - p} r_{t_1}(\boldsymbol{\Phi})^2 \left( 1 + \mathcal{O}\left( \frac{1}{t_2 - t_1 - p} \right) \right) + \frac{1}{t_2 - t_1 - p} \chi^2(\boldsymbol{\Phi}) \left( 1 + \mathcal{O}\left( \frac{1}{t_2 - t_1 - p} \right) \right) - \lambda . \tag{5}
\end{aligned}
$$

Neglecting the small terms $\mathcal{O}\left( \frac{1}{t_2 - t_1 - p} \right)$ leads to

$$\chi_\lambda^2(\boldsymbol{\Phi}) = r_{t_1}(\boldsymbol{\Phi})^2 - \lambda(t_2 - t_1 - p) . \tag{6}$$

Expression (6) has the following implications. Consider the case where all squared terms $r_i(\boldsymbol{\Phi})^2$ in the sum (1) defining $\chi_\lambda^2(\boldsymbol{\Phi})$ are approximately the same and independent of $t_1$, which occurs when the residuals are thin-tailed distributed and the model is well specified. Then, we have

$$r_i(\boldsymbol{\Phi})^2 \approx r^2 , \quad \forall i , \text{ including } r_{t_1}(\boldsymbol{\Phi})^2 = r^2 , \tag{7}$$

and thus

$$\chi_{np}^2(\boldsymbol{\Phi}) \approx r^2 . \tag{8}$$

Expressing (6) with the estimation (7) yields

$$\chi_\lambda^2(\boldsymbol{\Phi}) \approx r^2 - \lambda(t_2 - t_1 - p) . \tag{9}$$

Comparing with (8), this suggests that varying $t_1$ is expected in general to introduce a linear bias of the normalised sum $\chi_{np}^2(\boldsymbol{\Phi})$ of squares of the residuals, which is proportional to the size of the calibration window (up to the small correction by the number $p$ of degrees of freedom of the model). If we want to compare calibrations over different window sizes, we need to correct for this bias.

More specifically, rather than fixing the window size $t_2 - t_1 = N$, we want to determine the 'best' $t_1$, thus comparing calibrations for varying window sizes, for a fixed right end point $t_2$. As a consequence, the Lagrange multiplier $\lambda$ is no more fixed to ensure that the constraint $t_2 - t_1 = N$ holds, but now quantifies the average bias or "cost" associated with changing the window sizes. This bias is appreciable for small data sample sizes. It vanishes asymptotically as $N \to \infty$, i.e. $\lim_{N \to \infty} \lambda = 0$.

In Statistical Physics, this is analogous to the change from the canonical to the grand canonical ensemble, where the condition of a fixed number of particles (fixed number of points in a fixed window size) is relaxed to a varying number of particles with an energy cost per particle determined by the chemical potential (the Lagrange parameter $\lambda$) [14,15]. It is well-known that the canonical ensemble is recovered from the grand canonical ensemble by fixing the chemical potential (Lagrange multiplier) so that the number of particles is equal to the imposed constraints. Idem here.

How to determine the crucial Lagrange parameter $\lambda$? We propose an empirical heuristic. When plotting $\chi_{np}^2(\boldsymbol{\Phi})$ as a function of $t_1$ for various instances, we observe that a linearly decreasing function of $t_1$ provides a good approximation of it, as predicted by (6) (for $\lambda > 0$). The slope can then be interpreted as quantifying the average bias of the scaled goodness-of-fit $\chi_{np}^2(\boldsymbol{\Phi})$ due to the reduced number of data points as $t_1$ is increased (and the number $t_2 - t_1$ of points to fit

decreases). This average bias is clearly dependent on the data and of the model used to calibrate it. We can thus interpret the average linear trend observed empirically as determining the effective Lagrange regulariser term $\lambda$ that quantifies the impact on the goodness-of-fit resulting from the addition of data points in the calibration, given the specific realisation of the data and the model to calibrate. Thus, to make all the calibrations performed for different $t_1$ comparable for the determination of the optimal window size, we propose to correct expression (1) by subtracting the term $\lambda(t_2 - t_1)$ from the normalised sum of squared residuals $\chi^2_{np}(\boldsymbol{\Phi})$ given by Eq. (1), where $\lambda$ is estimated empirically as the large scale linear trend. Here, we omit the $p$ correction since it leads to a constant translation for a given model with given number of degrees of freedom. Such a large scale linear trend of $\chi^2_{np}(\boldsymbol{\Phi})$ as a function of $t_1$ has been reported for a number of financial bubble calibrations in Demos and Sornette [7]. Our proposed procedure thus amounts simply to detrend $\chi^2_{np}(\boldsymbol{\Phi})$, which has the effect of making more pronounced the minima of $\chi^2_{np}(\boldsymbol{\Phi})$, as we shall see below for different models.

To summarise, endogenising $t_1$ in the set of parameters to calibrate requires to minimise

$$\chi^2_\lambda(\boldsymbol{\Phi}) = \chi^2_{np}(\boldsymbol{\Phi}) - \lambda(t_2 - t_1) \tag{10}$$

$$= \frac{1}{(t_2 - t_1) - p} \sum_{i=t_1}^{t_2} r_i(\boldsymbol{\Phi})^2 - \lambda(t_2 - t_1) , \tag{11}$$

with

$$r_i(\boldsymbol{\Phi}) = y_i^{data} - y_i^{model}(\boldsymbol{\Phi}) , \tag{12}$$

where $\lambda$ is determined empirically so that $\chi^2_{np}(\boldsymbol{\Phi}) - \lambda(t_2 - t_1)$ has zero drift as a function of $t_1$ over the set of scanned values. The obtained empirical value of $\lambda$ can be used as a diagnostic parameter quantifying the tendency of the model to over-fit the data. We can thus also refer to $\lambda$ as the "overfit measure". When it is large, the goodness-of-fit $\chi^2(\boldsymbol{\Phi})$ changes a lot with the number of data points, indicating a poor overall ability of the model to account for the data. Demos and Sornette [7] observed other cases where $\chi^2(\boldsymbol{\Phi})$ is constant as a function of $t_1$ (corresponding to a vanishing $\lambda$), which can be interpreted in a regime where the model fits robustly the data, "synchronising" on its characteristic features in a way mostly independent of the number of data points.

In the next Section, we make use of a simulation study in order to precisely exemplify how $\chi^2_\lambda(\boldsymbol{\Phi})$ can be used to detect a regime change in the simple Linear Regression problem, thus yielding the optimal window length one should consider for fitting purposes.

## 3. Application of the Lagrange regularisation method to a simple linear-regression problem

### 3.1. Definitions and formal solution

Consider the following linear model

$$\boldsymbol{Y} = \beta\boldsymbol{X} + \boldsymbol{\varepsilon}, \tag{13}$$

with explanatory variable of length $(N \times 1)$ denoted by $\boldsymbol{X} = \{x_1, x_2, \ldots, x_N\}$, regressand $\boldsymbol{Y} = \{y_1, y_2, \ldots, y_N\}$ and error vector $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2)$. Bold variables denote either matrices or vectors. Fitting Eq. (13) to a given data set $Y^{data}$ consists on solving the quadratic minimisation problem

$$\hat{\beta} = \arg\min_\beta \chi^2, \tag{14}$$

where the Residual Sum of Squares $:= \chi^2$ reads

$$\chi^2 = \sum_{i=1}^N (Y_i^{data} - Y_i)^2. \tag{15}$$

For a given data set of length $N$, $\hat{\beta}$ is obtained via

$$\hat{\beta} = (\boldsymbol{X'X})^{-1}\boldsymbol{X'Y}. \tag{16}$$

Let $w^* \subseteq Y^{data}$ and have length $\leq$ N. $w^* \in [\tau : \bar{t}_2]$ thus denotes the optimal window size one should use for fitting a model into a data set of length N for a fixed end point $:= \bar{t}_2$ and an optimal starting point $:= \tau$. In the next section, we exemplify how $w^*$ can be obtained via the proposed Lagrange Regularisation technique through a simulation study.

### 3.2. Assessing the performance of the Lagrange regularisation method via a simulation study

In order to show how the goodness-of-fit metric $\chi^2(\boldsymbol{\Phi})$ fails to flag the optimal $\tau$-portion of the data set where the regime of interest exists and how delicate is $\chi^2_{np}(\boldsymbol{\Phi})$ for diagnosing the true value of the transition time $\tau$, 20000 synthetic realisations of the process (13) were generated, with $X := t \in [-200, +1]$, in such a way that $Y^{data}$ displays a sudden
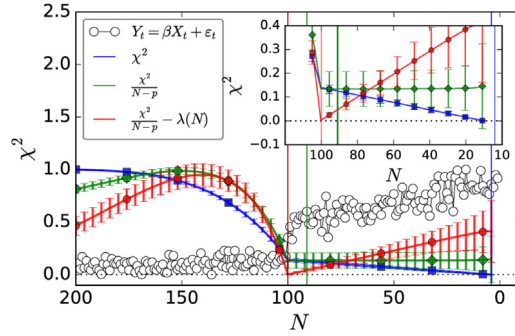
**Fig. 1.** Different goodness-of-fit measures applied to a shrinking-window linear regression problem (Eq. (13)) in order to diagnose the optimal calibration window length: We simulated synthetic time-series with length $N = 200$ (white circles) using expression (13) with a sudden change of regime at $t = -100$. We then fitted the same model (13) within shrinking windows (from left to right), i.e. for a fixed $t_2 = 1$, we shrink $t_1$ from $t_1 = -200$ to $t_1 = -3$ and show the values of $\chi^2(\Phi)$ (blue), $\chi_{np}^2(\Phi)$ (green) and $\chi_\lambda^2(\Phi)$ (red) metrics as a function of this shrinking estimation window. For each pair $[t_2 : t_1]$ (i.e. for each N), the process of generating synthetic data and fitting the model was repeated 20000 times (resulting on confidence bounds for each metric). For $t = [-200:-100]$, $Y_t$ was simulated with $\beta = 0.3$ while from $t = [-100:1]$, $\beta = 0.6$ was used. Without loss of generality, both the data and the cost functions had their values divided by their respectively maximum value in order to be bounded within the interval $[0, 1]$. A Python script for generating the figure and performing all calculations can be found in the Appendix.

change of regime at $\tau = -100$. In the first half of the data set $[-200, -100]$, the data points are generated with $\beta = 0.3$. In the second half of the data set $[-101, 0]$, the data points are generated with $\beta = 0.6$. As a consequence, the change includes both a jump and a doubling of slope, as shown in Fig. 1. After the addition of random noise $\epsilon \sim \mathcal{N}(0, 1)$, each single resulting time-series was fitted for a fixed end time $\bar{t}_2 = 1$ while shrinking the left-most portion of the data $(t_1)$ towards $t_2$, starting at $t_1 = -200, -199, \ldots, t_2 - 3$. For the largest window with $t_1 = -200$, there are $t_2 - t_1 + 1 = 1 - (-200) + 1 = 202$ data points to fit. For the smallest window with $t_1 = t_2 - 3$, there are $t_2 - t_1 + 1 = 4$ data points to fit. For each window size $w$, the process of generating synthetic data and fitting the model was repeated 20000 times, allowing us to obtain confidence intervals.

As depicted by Fig. 1, the proposed methodology is able to correctly diagnose the optimal starting point $:= \tau$ associated with the change of slope. While the $\chi^2(\Phi)$ metric monotonously decreases and the $\chi_{np}^2(\Phi)$ metric plateaus from $t = -100$ onwards, $\chi_{np}^2 - \lambda(t_2 - t_1)$ monotonously increases over the same interval, thus marking a clear minimum. The variance of the metric $\chi_\lambda^2(\Phi)$ also increases over this interval. Specifically, the metric $\chi^2(\Phi)$ tends to favour the smallest windows and therefore overfitting is prone to develop and remain undetected. The metric $\chi_{np}^2(\Phi)$ suggests $\tau \approx -90$ after 20000 simulations, which is 10% away from the true value $\tau = 100$. Moreover, the dependence of $\chi_{np}^2(\Phi)$ as a function of $t_1$ is so flat for $t_1 \in [-100 : -40]$ that any given value of $\tau$ within this period is statistically significant. For this simulation study, $\chi_{np}^2(\Phi)$ ranges for 0.134 to 0.135 for $t_1 \in [-100 : -60]$, so as to be almost undistinguishable over this interval of possible $\tau$ values. As we shall see later on, the performance of $\chi_{np}^2(\Phi)$ degrades further to resemble that of the $\chi^2(\Phi)$ metric when dealing with more complex nonlinear models such as the LPPLS model. On the other hand, our proposed correction via the Lagrange regulariser $\lambda$ provides a simple and effective method to identify the change of regime and the largest window size compatible with the second regime. The minimum is very pronounced and clear, which is not the case for $\chi_{np}^2(\Phi)$.

## 4. Using the Lagrange regularisation method for detecting the beginning of financial bubbles

In the previous section, we have tested the proposed novel goodness-of-fit metric for inferring the optimal beginning point or change point $\tau$ (for a fixed end point $\bar{t}_2$) in the calibration of a simple linear model. The application of the Lagrange regulariser $\lambda$ allowed us to find the optimal window length $w^* = [\tau : t_2]$ for fitting the model by enabling the comparison of the goodness-of-fits across different $w$ values. We now extend the application of the methodology to a more complex non-linear model, which requires to compare fits across different window sizes in order to diagnose bubble periods on financial instruments such as equity prices and price indexes.

### 4.1. The LPPLS model and its calibration

The LPPLS (log-periodic power law singularity) model [8,16] provides a flexible set-up for diagnosing periods of price exuberance [17] on financial instruments. It highlights the role of herding behaviour, translating into positive feedbacks in the price dynamics during the formation of bubbles. This is reflected in faster-than-exponential growth of the price of financial instruments. Such explosive behaviour is completely unsustainable and the bubbles usually ends with a crash or a progressive correction. Here, we use the LPPLS model combined with the Lagrange regulariser $\lambda$ in order to detect the beginning of financial bubbles.

In the LPPLS model, the expectation of the logarithm of the price of an asset is written under the form

$$fLPPL(\boldsymbol{\phi}, t) = A + B(f) + C_1(g) + C_2(h), \tag{17}$$

where $\boldsymbol{\phi} = \{A, B, C_1, C_2, m, \omega, t_c\}$ is a $(1 \times 7)$ vector of parameters we want to determine and

$$f \equiv (t_c - t)^m, \tag{18}$$

$$g \equiv (t_c - t)^m \cos(\omega \ln(t_c - t)), \tag{19}$$

$$h \equiv (t_c - t)^m \sin(\omega \ln(tc - t)). \tag{20}$$

Note that the power law singularity $(t_c - t)^m$ embodies the faster-than-exponential growth. Log-periodic oscillations represented by the cosine and sine of $\ln(t_c - t)$ model the long-term volatility dynamics decorating the accelerating price. Expression (17) uses the formulation of Filimonov and Sornette [18] in terms of 4 linear parameters $A, B, C_1, C_2$ and 3 nonlinear parameter $m, \omega, t_c$.

Fitting Eq. (17) to the log-price time-series amounts to search for the parameter set $\boldsymbol{\phi}^*$ that yields the smallest $N$-*dimensional* distance between realisation and theory. Mathematically, using the $L^2$ norm, we form the following sum of squares of residuals

$$F(t_c, m, \omega, A, B, C_1, C_2) = \sum_{i=1}^{N} \left[ \ln[P(t_i)] - A - B(f_i) - C_1(g_i) - C_2(h_i) \right]^2, \tag{21}$$

for $i = 1, \ldots, N$. We proceed in two steps. First, enslaving the linear parameters $\{A, B, C_1, C_2\}$ to the remaining nonlinear parameters $\boldsymbol{\phi} = \{t_c, m, \omega\}$, yields the cost function $\chi^2(\boldsymbol{\phi})$

$$\chi^2(\boldsymbol{\phi}) := F_1(t_c, m, \omega) \tag{22}$$

$$= \min_{\{A,B,C_1,C_2\}} F(t_c, m, \omega, A, B, C_1, C_2) \tag{23}$$

$$= F(t_c, m, \omega, \widehat{A}, \widehat{B}, \widehat{C}_1, \widehat{C}_2), \tag{24}$$

where the hat symbol ⌢ indicates estimated parameters. This is obtained by solving the optimisation problem

$$\{\widehat{A}, \widehat{B}, \widehat{C}_1, \widehat{C}_2\} = arg \min_{\{A,B,C_1,C_2\}} F(t_c, m, \omega.A, B, C_1, C_2), \tag{25}$$

for fixed $t_c, m$ and $\omega$, which can be obtained analytically by solving the following system of equations,

$$\begin{bmatrix} N & \sum f_i & \sum g_i & \sum h_i \\ \sum f_i & \sum f_i^2 & \sum f_i g_i & \sum f_i h_i \\ \sum g_i & \sum f_i g_i & \sum g_i^2 & \sum g_i h_i \\ \sum h_i & \sum f_i h_i & \sum g_i h_i & \sum h_i^2 \end{bmatrix} \begin{bmatrix} \widehat{A} \\ \widehat{B} \\ \widehat{C}_1 \\ \widehat{C}_2 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum y_i f_i \\ \sum y_i g_i \\ \sum y_i h_i \end{bmatrix}. \tag{26}$$

Second, we solve the nonlinear optimisation problem involving the remaining nonlinear parameters $m, \omega, t_c$:

$$\{\widehat{t_c}, \widehat{m}, \widehat{\omega}\} = arg \min_{\{t_c,m,\omega\}} F_1(t_c, m, \omega). \tag{27}$$

The model is calibrated on the data using the Ordinary Least Squares method, providing estimations of all parameters $t_c$, $\omega$, $m$, $A$, $B$, $C_1$, $C_2$ in a given time window of analysis.

For each fixed data point $t_2$ (corresponding to a fictitious "present" up to which the data is recorded), we fit the price time series in shrinking windows $(t_1, t_2)$ of length $dt := t_2 - t_1$ decreasing from 1600 trading days to 30 trading days. We shift the start date $t_1$ in steps of 3 trading days, thus giving us 514 windows to analyse for each $t_2$. In order to minimise calibration problems and address the sloppiness of the model with respect to some of its parameters (and in particular $t_c$), we use a number of filters to select the solutions. For further information about the sloppiness of the LPPLS model, we refer to [7,12,19,20]. The filters used here are $\{(0.1 < m < 0.9), (6 < \omega < 13), (t_2 - [t2 - t1] < t_c < t_2 + [t2 - t1])\}$, so that only those calibrations that meet these conditions are considered valid and the others are discarded. These filters derive from the empirical evidence gathered in investigations of previous bubbles [20–22].

Previous calibrations of the JLS model have further shown the value of additional constraints imposed on the nonlinear parameters in order to remove spurious calibrations (false positive identification of bubbles) [7,23,24]. For our purposes, we do not consider them here.

### 4.2. Empirical analysis

We apply our novel goodness-of-fit metric to the problem of finding the beginning times of financial bubbles, defined as the optimal starting time $t_1$ obtained by endogenising $t_1$ and calibrating it. We first illustrate and test the method on synthetic time series and then apply it to real-world financial bubbles. A Python implementation of the algorithm is provided in the Appendix.

#### 4.2.1. Construction of synthetic LPPLS bubbles and results

To gain insight about the application of our proposed calibration methodology on a controlled framework and thus establish a solid background to our empirical analysis, we generate synthetic price time series that mimic the salient properties of financial bubbles, namely, a power law-like acceleration decorated by oscillations. The synthetic price time series are obtained by using formula (17) with parameters given by the best LPPLS fit within the window $w \in [t_1 = 1 \, Jan. \, 1981: t_2 = 30 \, Aug. \, 1987]$ of the bubble that ended with the Black Monday 19 Oct. 1987 crash. These parameters are m = 0.44, $\omega$=6.5, $C_1 = -0.0001$, $C_2 = 0.0005$, A = 1.8259, B = $-0.0094$, $t_c = 1194$ (corresponding to 1987/11/14), where days are counted since an origin put at $t_1 = Jan. \, 1981$. To the deterministic component describing the expected log-price given by expression (17) and denoted by $fLPPLS(\phi, t)$, we add a stochastic element to obtain the synthetic price time series

$$\ln[P(t)] = fLPPLS(\phi, t) + \sigma \epsilon(t), \tag{28}$$

where $\epsilon(t) \sim \mathcal{N}(0, \sigma^2)$ noise, $\sigma^2 = 0.03$ and $t = [1, \ldots, N = 1100]$.

To illustrate how we create a price time series with a well-defined transition point corresponding to the beginning of a bubble, we take the first 500 points generated with expression (28) and mirror them via a $t \to t_1 - t$ reflection across the time period: {1 Jan. 1910 : 1 jun. 1911}. In turn , we concatenate this reflected sequence of 500 prices to the 1100 prices obtained with (28) for $t \geq t_1$, so that the true transition point corresponding to the start of the bubble described by the LPPLS pattern is $t_1 = 1 \, jun. \, 1911$. The black stochastic line on the top of Fig. 2 represent this union of the two time-series. This union constitutes the whole synthetic time series on which we are going to apply our Lagrange regularisation of $\chi^2_\lambda(\Phi)$ in order to attempt recovering the true start time, denoted by the hypothetical time $t_1 = 1 \, jun. \, 1911$. The same procedure is applied to the other real financial time series.

For the synthetic bubble price time series shown in the top panel of Fig. 2, we thus calibrated it with Eq. (17) by minimising expression (1) in windows $w = [t_1, t_2]$, varying $t_2$ from 1912/07/01 to $t_2 = 1913/01/01$, with $t_1$ scanned from $t_1 = Jan. \, 1910$ up to 30 business days before $t_2$, i.e. up to $t_{1,max} = t_2 - 30$ for each fixed $t_2$. The goal is to determine whether the transition point $\tau$ we determine is close (or even equal to) the true hypothetical value $t_1 = 01 \, jun. \, 1911$ for different maturation times $t_2$ of the bubble. The number of degrees of freedom used for this exercise as well as for the real-world time series is $p = 8$, which includes the 7 parameters of the LPPLS model augmented by the extra parameter $t_1$.

The results of the determination of the beginning times of the synthetic time-series with end time $t_2 = \{1912.07.01; \, 1912.10.01; \, 1912.11.15; \, 1913.01.01\}$ are shown in Fig. 3. For the earliest $t_2$ = 1912/07/01, our proposed goodness-of-fit scheme is already capable of roughly diagnosing correctly the bubble beginning time, finding the optimal $\tau$ to be $\approx May \, 1911$. In contrast, the competing metric $(\chi^2_{np}(\Phi))$ is degenerate as $t_1 \to t_2$ and is thus blind to the beginning of the bubble. For $t_2$ closer to the end of the bubble, $\chi^2_{np}(\Phi)$ continues to deliver very small optimal windows, proposing the incorrect conclusion that the bubble has started very recently (i.e. close to the pseudo present time $t_2$). This is a signature of strong overfitting, which is quantified via $\lambda$ and depicted in the title of the figure alongside with the bubble beginning time and $t_2$. The Lagrange regularisation of the $\chi^2_\lambda(\Phi)$ locks into the true value of $\tau \approx jul.1911$ as $t_2 \to t_c$, i.e., as $t_2$ moves closer and closer to January 1913 and the LPPLS signal becomes stronger. In summary, Fig. 3 shows that the optimal starting time $\tau$ of the bubble determined by using $\chi^2_{np}(\Phi)$ is very close to its $t_2$. In contrast, the optimal $\tau$ for $\chi^2_\lambda(\Phi)$ is close to the true beginning of the LPPLS regime. This confirms the excellent performance of our proposed method.

#### 4.2.2. Real-world data: set-up of the analysis of bubble periods of different financial indices

The real-world data sets analysed here consist in bubble periods that have occurred on the following major Indexes: S&P-500 with the list of $t_2$'s given here,[1] IBovespa with the list of $t_2$'s given here[2] and SSEC with the list of $t_2$'s given here.[3] For each data set and for each fixed pseudo present time $t_2$ depicted by red vertical dashed lines in Fig. 2, our search for the bubble beginning time $\tau$ consists in fitting the LPPLS model using a shrinking estimation window $w$ with $t_1 = [t_2 - 30 : t_2 - 1600]$ with incremental step-size of 3 business days. This yields a total of 514 fits per $t_2$.

#### 4.2.3. Results of the determination of the beginning times of real-world bubbles

For the S&P-500 Index, the results shown in Fig. 4 are even more pronounced than for the synthetic time series. While again $\chi^2_{np}(\Phi)$ is unable to diagnose the optimal starting date of a faster than exponential log-price growth $\tau \equiv t_1$, the Lagrange regularisation of the $\chi^2_{np}(\Phi)$ depicted by blank triangles in the lower box of the figure is capable of overcoming the tendency of the model to overfit data as $t_1 \to t_2$. Specifically, the method diagnoses the start of the Black-Monday bubble at $t_1 \approx March \, 1984$ and the beginning of the Sub-Prime bubble at $\approx Aug. \, 2003$ in accordance with [25].

We also picked two pseudo present times $t_2's$ at random in order to check how consistent are the results. To our delight, the method is found capable of capturing the different time-scales present in the bubble formation in an endogenous

---

[1] $t_2$'s = {1987.07.15; 1997.06.01; 2000.01.01; 2007.06.01}.

[2] $t_2$'s = {2000.01.01; 2004.01.01; 2006.01.01; 2007.12.01}.

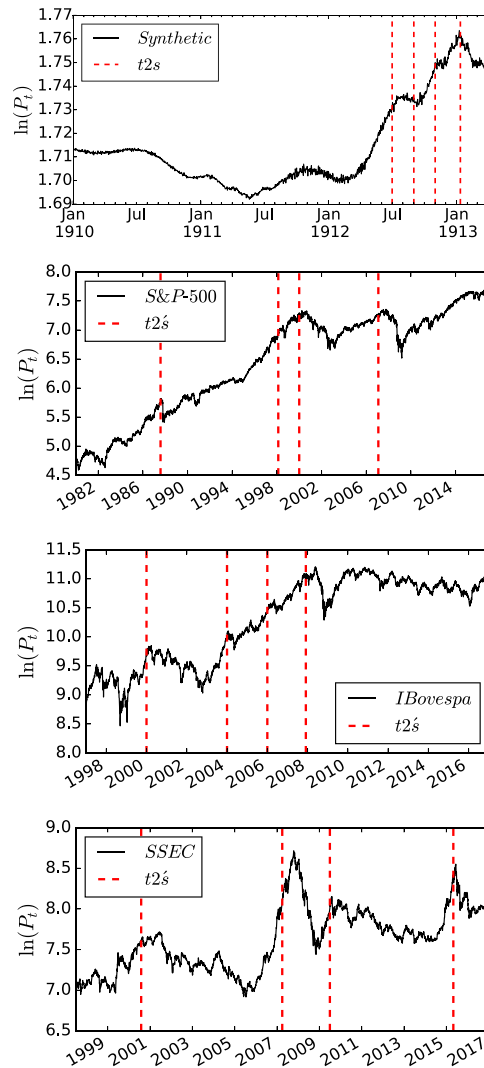[3] $t_2$'s = {2000.08.01; 2007.05.01; 2009.07.01; 2015.05.01}.

**Fig. 2.** Synthetic and real-world Time-series used in this study for measuring the performance of different goodness-of-fit metrics at different $t_2$'s (red lines): Synthetic time-series and Indexes S&P-500, *IBovespa* and *SSEC* with the values of $t_2$'s used in the subsequent figures given by $t_2's$ = {1912.07.01; 1912.10.01; 1912.11.15; 1913.01.01}, $t_2's$ = {1987.07.15; 1997.06.01; 2000.01.01; 2007.06.01}, $t_2's$ = {2000.01.01; 2004.01.01; 2006.01.01; 2007.12.01} and $t_2's$ = {2000.08.01; 2007.05.01; 2009.07.01; 2015.05.01} respectively (red dashed vertical lines).

manner. For $t_2 = 1997.06.01$, the method suggests the presence of a bubble that nucleated more than five years earlier. This recovers the bubble and change of regime in September 1992, documented in Chapter 9 of Sornette [26] as a genuine change of regime since the market stopped its ascent and plateaued for the three following months. For $t_2 = 2000.01.01$, $\chi^2_\lambda(\Phi)$ diagnoses a bubble with a shorter duration, which started in November 1998. The starting time is coherent with the recovery after the so-called Russian crisis of August–September 1998 when the US stock markets dropped by about 20%. And this bubble is nothing but the echo in the S&P500 of the huge dotcom bubble that crashed in March–April 2000. More generally, scanning $t_2$ and different intervals for $t_1$, the Lagrange regularisation of the $\chi^2_{np}(\Phi)$ can endogenously identify a hierarchy of bubbles of different time-scales, reflecting their multi-scale structure [12,26].

For the IBovespa and the SSEC Index (Figs. 5 and 6 respectively), the huge superiority of the Lagrange regularisation of the $\chi^2_{np}(\Phi)$ vs. the $\chi^2_{np}(\Phi)$ metric is again obvious. For each of the four chosen $t_2$'s in each figure, $\chi^2_\lambda(\Phi)$ exhibits a well-marked minimum corresponding to a well-defined starting time for the corresponding bubble. These objectively identified $t_1$ correspond pleasantly to what the eye would have chosen. They pass the "smell test" [27]. In contrast, the $\chi^2_{np}(\Phi)$ metric provides essentially no guidance on the determination of $t_1$.
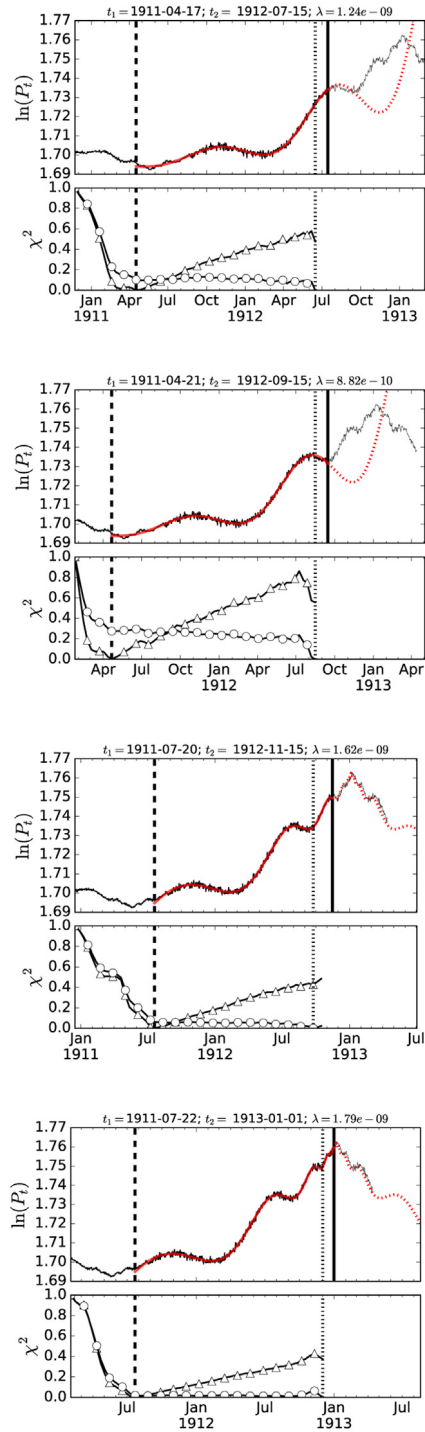
**Fig. 3.** Diagnosing the beginning of financial bubbles by comparing two goodness-of-fit metrics $\chi^2_{np}(\Phi)$ vs. $\chi^2_{\lambda}(\Phi)$ using the LPPLS model on Synthetic Time-Series: in the lower plot of each $t_2$ value, $\chi^2_{np}(\Phi)$ is depicted as empty circles while $\chi^2_{\lambda}(\Phi)$ is shown as empty triangles. The dashed thick black vertical lines (resp. thin hashed vertical lines) show the position of the minimum of each goodness-of-fit metric and therefore represents the optimal $\tau \equiv t_1$ for $\chi^2_{\lambda}(\Phi)$ and $\chi^2_{np}(\Phi)$ respectively. Note that the optimal $\tau$ for $\chi^2_{np}(\Phi)$ is very close to its $t_2$. In contrast, the optimal $\tau$ for $\chi^2_{\lambda}(\Phi)$ is close to the true beginning of the LPPLS regime. For a fixed $t_2$, the log-price time-series of the Index was fitted using a shrinking window from $t_1 = [t_2 - 30 : t_2 - 1600]$ sampled every 3 days. For a fixed $t_2$ and $t_1$, we display the resulting fit of the LPPLS model (red line) obtained with the parameters solving Eq. (27).
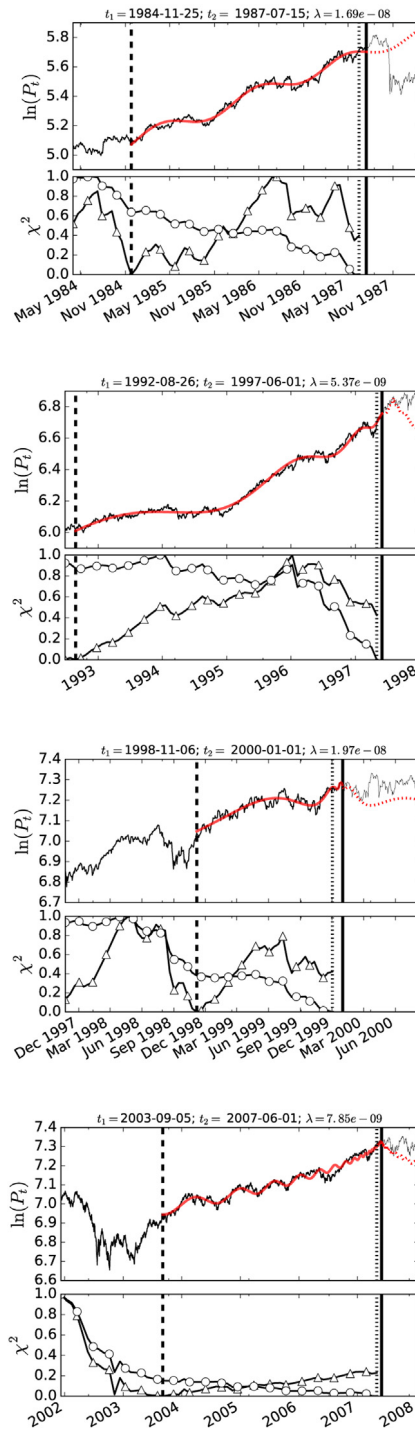
**Fig. 4.** Same as Fig. 3 for the US S&P-500 Index.

## 5. Conclusion

We have presented a novel goodness-of-fit metric, aimed at comparing goodnesses-of-fit across a nested hierarchy of data sets of shrinking sizes. This is motivated by the question of identifying the start time of financial bubbles, but applies more generally to any calibration of time series in which the start time of the latest regime of interest is unknown. We have introduced a simple and physically motivated way to correct for the overfitting bias associated with shrinking data
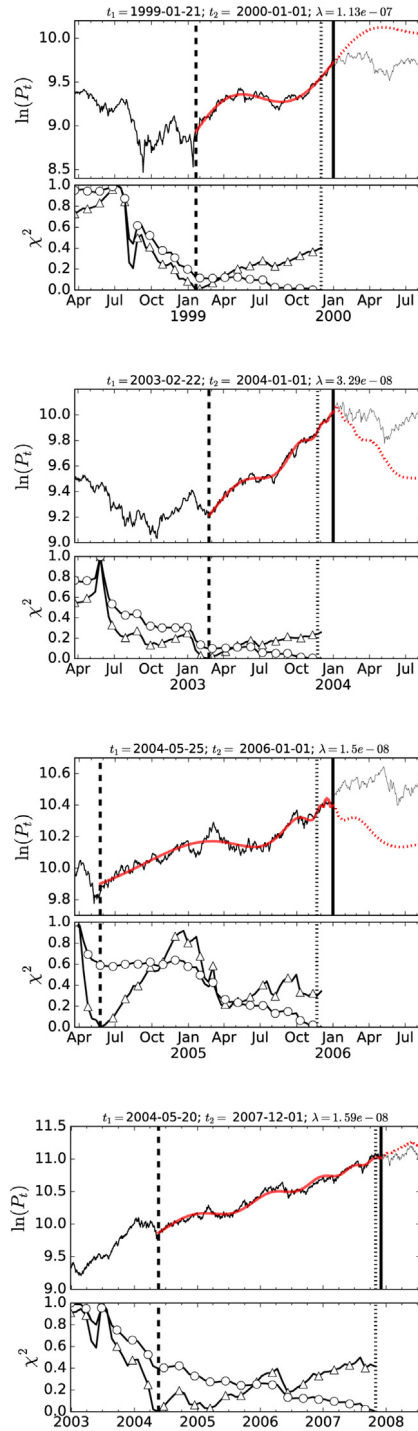
**Fig. 5.** Same as Fig. 3 for the Brazilian *IBovespa* index.

sets, which we refer to at the Lagrange regularisation of the $\chi^2_{np}(\boldsymbol{\Phi}) := \frac{1}{N-p} SSR$. We have suggested that the bias can be captured by a Lagrange regularisation parameter $\lambda$. In addition to helping remove or alleviate the bias, this parameter can be used as a diagnostic parameter, or "overfit measure", quantifying the tendency of the model to overfit the data. It is a function of both the specific realisation of the data and of how the model matches the generating process of the data.
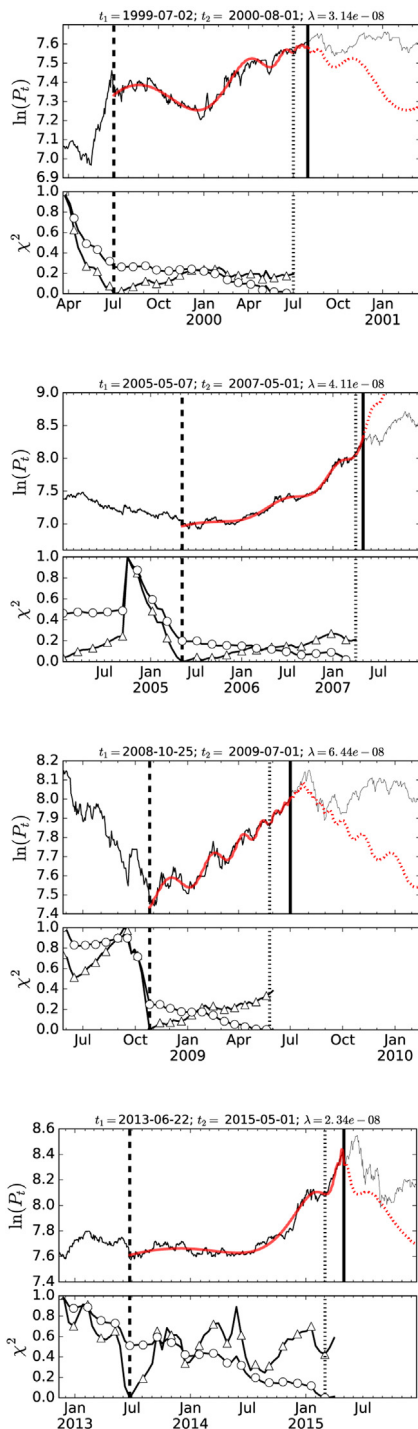
**Fig. 6.** Same as Fig. 3 for the Chinese *SSEC* index.

Applying the Lagrange regularisation of the $\chi^2_{np}(\Phi)$ to simple linear regressions with a change point, synthetic models of financial bubbles with a well-defined transition regime and to a number of financial time series (US S&P500, Brazil IBovespa and China SSEC Indices), we document its impressive superiority compared with the $\chi^2_{np}(\Phi)$ metric. In absolute sense, the Lagrange regularisation of the $\chi^2_{np}(\Phi)$ is found to provide very reasonable and well-defined determinations of

the starting times for major bubbles such as the bubbles ending with the 1987 Black-Monday, the 2008 Sub-prime crisis and minor speculative bubbles on other Indexes, without any further exogenous information.

Overall, the usage of the proposed technique for correcting the overfitting tendency of a model $\chi^2_\lambda(\Phi)$, alongside with the LPPLS model allows one to answer in real-time an extremely important question in finance, namely, the question of when a period of price exuberance (i.e. financial bubble) started. The elegant solution proposed here can be straightforwardly implemented such that, conditional on a pseudo-current time, the beginning time of a bubble ($t_1^*$) can be endogenously obtained without any further exogenous information requirements. Moreover, the usage of the technique goes beyond the specific scope of finance, being essentially implementable on any statistical problem whose cost function is degenerate as the sample size shrinks as a function of time. Examples of such problems consists on: defining the optimal window length for computing a moving average, comparing goodness-of-fit across unbalanced sample sizes for a given model, automatic change point detection, among others.

## Appendix

```python
- Python script for computing the Lambda regulariser metric: OLS example.
- Copyright: G.Demos @ ETH-Zurich.
- Jun.2018

#----------------------------
def simulateOLS():
    """ Generate synthetic OLS as presented in the paper """
    nobs = 200
    X    = np.arange(0,nobs,1)
    e    = np.random.normal(0, 10, nobs)
    beta = 0.5
    Y    = [beta*X[i] + e[i] for i in range(len(X))]
    Y = np.array(Y)
    X = np.array(X)
    Y[:100] = Y[:100] + 4*e[:100]
    Y[100:200] = Y[100:200]*8
    return X, Y


#----------------------------
def fitDataViaOlsGetBetaAndLine(X,Y):
    """ Fit synthetic OLS """
    beta_hat = np.dot(X.T,X)**-1. * np.dot(X.T,Y) # get beta
    Y = [beta_hat*X[i] for i in range(len(X))] # generate fit
    return Y


#----------------------------
def getSSE(Y, Yhat, p=1, normed=False):
    """
    Obtain SSE (chi^2)
    p -> No. of parameters
    Y -> Data
    Yhat -> Model
    """
    error = (Y-Yhat)**2.
    obj = np.sum(error)
    if normed == False:
        obj = np.sum(error)
    else:
        obj = 1/np.float(len(Y) - p) * np.sum(error)
    return obj


#----------------------------
def getSSE_and_SSEN_as_a_func_of_dt(normed=False, plot=False):
    """ Obtain SSE and SSE/N for a given shrinking fitting window w """
    # Simulate Initial Data
    X, Y = simulateOLS()
```

```python
    # Get a piece of it: Shrinking Window
    _sse = []
    _ssen = []
    for i in range(len(X)-10): # loop t1 until: t1 = (t2 - 10):
        xBatch = X[i:-1]
        yBatch = Y[i:-1]
        YhatBatch = fitDataViaOlsGetBetaAndLine(xBatch, yBatch)
        sse = getSSE(yBatch, YhatBatch, normed=False)
        sseN = getSSE(yBatch, YhatBatch, normed=True)
        _sse.append(sse)
        _ssen.append(sseN)

    if plot == False:
        pass
    else:
        f, ax = plt.subplots(1,1,figsize=(6,3))
        ax.plot(_sse, color='k')
        a = ax.twinx()
        a.plot(_ssen, color='b')
        plt.tight_layout()
    if normed==False:
        return _sse, _ssen, X, Y # returns results + data
    else:
        return _sse/max(_sse), _ssen/max(_ssen), X, Y # returns results + data


#-----------------------------
def LagrangeMethod(sse):
    """ Obtain the Lagrange regulariser for a given SSE/N"""
    # Fit the decreasing trend of the cost function
    slope = calculate_slope_of_normed_cost(sse)

    return slope[0]


#-----------------------------
def calculate_slope_of_normed_cost(sse):
    #Create linear regression object using statsmodels package
    regr = linear_model.LinearRegression()

    # create x range for the sse_ds
    x_sse = np.arange(len(sse))
    x_sse = x_sse.reshape(len(sse),1)

    # Train the model using the training sets
    res = regr.fit(x_sse, sse)
    return res.coef_


#-----------------------------
def obtainLagrangeRegularizedNormedCost(X, Y, slope):
    """ Obtain the Lagrange regulariser for a given SSE/N Pt. III"""
    Yhat = fitDataViaOlsGetBetaAndLine(X,Y) # Get Model fit
    ssrn_reg = getSSE(Y, Yhat, normed=True) # Classical SSE
    ssrn_lgrn = ssrn_reg - slope*len(Y) # SSE lagrange
    return ssrn_lgrn


#-----------------------------
def GetSSEREGvectorForLagrangeMethod(X, Y, slope):

    """
    X and Y used for calculating the original SSEN
    slope is the beta of fitting OLS to the SSEN
    """
```

```
# Estimate the cost function pondered by lambda using a Shrinking Window.
_ssenReg = []
for i in range(len(X)-10):
    xBatch = X[i:-1]
    yBatch = Y[i:-1]
    regLag = obtainLagrangeRegularizedNormedCost(xBatch,
            yBatch,
            slope)
    _ssenReg.append(regLag)
return _ssenReg
```

## References

[1] S. Loscalzo, L. Yu, C. Ding, Consensus group stable feature selection, in: Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2009, pp. 567–576.
[2] R. Tibshirani, Regression shrinkage and selection via the lasso, R. Stat. Soc. 58 (1) (1996) 267–288.
[3] A.Y. Ng, Feature selection, l1 vs. l2 regularization, and rotational invariance, in: Proceedings of the Twenty-First International Conference on Machine Learning, ACM, 2004, p. 78,
[4] H. Akaike, A new look at the statistical model identification, IEEE Trans. Automat. Control 19 (6) (1974) 716–723.
[5] F. Peres, F. Fogliatto, Variable selection methods in multivariate statistical process control: a systematic literature review, Comput. Ind. Eng. 115 (2018) 603–619.
[6] S. Aminikhanghahi, D. Cook, A survey of methods for time series change point detection, Knowl. Inf. Syst. 51 (2017) 339–367.
[7] G. Demos, D. Sornette, Birth or burst of financial bubbles: which one is easier to diagnose?, Quant. Finance 5 (2017) 657–675.
[8] A. Johansen, O. Ledoit, D. Sornette, Crashes as critical points, Int. J. Theor. Appl. Finance 2 (2000) 219–255.
[9] D. Sornette, A. Johansen, Significance of log-periodic precursors to financial crashes, Quant. Finance 1 (4) (2001) 452–471.
[10] D. Sornette, W.-X. Zhou, Predictability of large future changes in major financial indices, Int. J. Forecast. 22 (2006) 153–168.
[11] V. Filimonov, G. Demos, D. Sornette, Financial bubbles: mechanisms, diagnostics and state of the world, Rev. Behav. Econ. 2 (3) (2015) 279–305.
[12] V. Filimonov, G. Demos, D. Sornette, Modified profile likelihood inference and interval forecast of the burst of financial bubbles, Quant. Finance 7 (8) (2017) 1167–1186.
[13] E. Fama, T. Shiller, L.P. Hansen, Understanding asset prices, The Royal Swedish Academy of Sciences (2013).
[14] J.W. Gibbs, Elementary Principles in Statistical Mechanics, in: Dover Books on Physics, Dover Publications, 1902.
[15] H.B. Callen, Thermodynamics and an Introduction to Thermostatistics (2nd ed.), New York: John Wiley & Sons, 1985.
[16] A. Johansen, D. Sornette, O. Ledoit, Predicting financial crashes using discrete scale invariance, J. Risk 1 (4) (1999) 5–32.
[17] R. Shiller, Irrational exuberance, Princeton University Press, Princeton, NJ, 2000.
[18] V. Filimonov, D. Sornette, A stable and robust calibration scheme of the log-periodic power law model, Physica A 392 (17) (2013) 3698–3707.
[19] D.S. Brée, D. Challet, P.P. Peirano, Prediction accuracy and sloppiness of log-periodic functions, Quant. Finance 13 (2) (2013) 275–280.
[20] D. Sornette, G. Demos, Q. Zhang, P. Cauwels, V. Filimonov, Q. Zhang, Real-time prediction and post-mortem analysis of the shanghai 2015 stock market bubble and crash, J. Invest. Strateg. 4 (4) (2015) 77–95.
[21] W.-X. Zhou, D. Sornette, Evidence of a worldwide stock market log-periodic anti-bubble since mid-2000, Physica A 330 (3–4) (2003) 543–583.
[22] Q. Zhang, Q. Zhang, D. Sornette, Early warning signals of financial crises with multi-scale quantile regressions of log-periodic power law singularities, PLoS ONE 11 (11) (2016) e0165819, 1–43, http://dx.doi.org/10.1371/journal.pone.0165819.
[23] D. Bree, D. Challet, P. Peirano, Prediction accuracy and sloppiness of log-periodic functions, Quant. Finance 3 (2013) 275–280.
[24] P. Geraskin, D. Fantazzini, Everything you always wanted to know about log-periodic power laws for bubble modeling but were afraid to ask, Eur. J. Finance 19 (5) (2011) 366–391.
[25] W.-X. Zhou, D. Sornette, Testing the stability of the 2000-2003 US stock market "antibubble", Physica A 348 (2005) 428–452.
[26] D. Sornette, Why stock markets crash: Critical events in complex financial systems, Princeton University Press, New Jersey, 2003.
[27] R. Solow, Building a science of economics for the real world, in: House Committee on Science and Technology; Subcommittee on Investigations and Oversight, 2010, July 20.