

PAPER • OPEN ACCESS

## Binary versus non-binary information in real time series: empirical results and maximum-entropy matrix models

To cite this article: Assaf Almog and Diego Garlaschelli 2014 *New J. Phys.* **16** 093015

View the [article online](#) for updates and enhancements.

### Related content

- [Analytical maximum-likelihood method to detect patterns in real networks](#)  
Tiziano Squartini and Diego Garlaschelli
- [On the concentration of large deviations for fat tailed distributions, with application to financial data](#)  
Mario Filiasi, Giacomo Livan, Matteo Marsili et al.
- [Topical Review](#)  
Andrea De Martino and Matteo Marsili

### Recent citations

- [Mesoscopic Community Structure of Financial Markets Revealed by Price and Sign Fluctuations](#)  
Assaf Almog *et al*

## Binary versus non-binary information in real time series: empirical results and maximum-entropy matrix models

**Assaf Almog and Diego Garlaschelli**

Instituut-Lorentz for Theoretical Physics, Leiden Institute of Physics, University of Leiden, Niels Bohrweg 2, 2333 CA Leiden, The Netherlands  
E-mail: [almog@lorentz.leidenuniv.nl](mailto:almog@lorentz.leidenuniv.nl)

Received 12 May 2014, revised 22 July 2014  
Accepted for publication 7 August 2014  
Published 12 September 2014

*New Journal of Physics* **16** (2014) 093015  
[doi:10.1088/1367-2630/16/9/093015](https://doi.org/10.1088/1367-2630/16/9/093015)

### Abstract

The dynamics of complex systems, from financial markets to the brain, can be monitored in terms of multiple time series of activity of the constituent units, such as stocks or neurons, respectively. While the main focus of time series analysis is on the magnitude of temporal increments, a significant piece of information is encoded into the binary projection (i.e. the sign) of such increments. In this paper we provide further evidence of this by showing strong nonlinear relations between binary and non-binary properties of financial time series. These relations are a novel quantification of the fact that extreme price increments occur more often when most stocks move in the same direction. We then introduce an information-theoretic approach to the analysis of the binary signature of single and multiple time series. Through the definition of maximum-entropy ensembles of binary matrices and their mapping to spin models in statistical physics, we quantify the information encoded into the simplest binary properties of real time series and identify the most informative property given a set of measurements. Our formalism is able to accurately replicate, and mathematically characterize, the observed binary/non-binary relations. We also obtain a phase diagram allowing us to identify, based only on the instantaneous aggregate return of a set of multiple time series, a regime where the so-called ‘market mode’ has an optimal interpretation in terms of collective (endogenous)



Content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](https://creativecommons.org/licenses/by/3.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

effects, a regime where it is parsimoniously explained by pure noise, and a regime where it can be regarded as a combination of endogenous and exogenous factors. Our approach allows us to connect spin models, simple stochastic processes, and ensembles of time series inferred from partial information.

Keywords: time series, econophysics, maximum entropy matrices

## 1. Introduction

In large systems, the observed dynamics or activity of each unit can be represented by a discrete time series providing a sequence of measurements of the state of that unit. One of the main challenges researchers are faced with is that of extracting meaningful information from the high-dimensional (multiple) time series characterizing all the elements of a complex system [1–9]. Traditionally, the main object of time series analysis is the characterization of patterns in the amplitude of the increments of the quantities of interest. Given a signal  $s_i(t)$  where  $i$  denotes one of the  $N$  units of the system and  $t$  denotes one of the  $T$  observed temporal snapshots, the generic increment or ‘return’  $r_i(t)$  can be defined as

$$r_i(t) \equiv s_i(t+1) - s_i(t) \quad i = 1, N \quad t = 1, T \quad (1)$$

and generates a new time series.

While a time series of increments encapsulates all the relevant information about the amplitude of the fluctuations of the original signal, a significant part of this information is encoded in the purely ‘binary’ projection of  $r_i(t)$ , i.e. its sign

$$x_i(t) \equiv \text{sign}[r_i(t)] = \begin{cases} +1 & r_i(t) > 0 \\ 0 & r_i(t) = 0 \\ -1 & r_i(t) < 0 \end{cases} \quad (2)$$

Previous analyses, mainly in the field of finance, have indeed documented various forms of statistical dependency between the sign and the absolute value of fluctuations, e.g. sign–volume correlations [10, 11] and the leverage effect [12–15]. Other studies have also documented that the binary projections of various financial [16] and neural [17] time series exhibit non-trivial dynamical features that resemble those of the original data. All these results suggest that binary projections indeed retain a non-trivial piece of information about the original time series, and call for a deeper analysis of the problem.

Being binary, the sign of the increments is much more robust to noise than the increments themselves. Moreover, it is scale-invariant (i.e. independent of the chosen unit of increments) and does not depend on whether the original data have been preliminarily rescaled or log-transformed (as usually done, e.g., for financial time series). Binary time series can also be analyzed with the aid of much simpler mathematical models than required by non-binary data (several examples of such models will be provided in this paper). Finally, as we show later on, in multiple financial time series the total binary increment of a given cross-section measures the instantaneous level of synchronization (i.e. the number of stocks moving in the same direction) of the market, while the total non-binary increment does not carry this piece of information. For all the above reasons, it is important to further investigate whether the full ‘weighted’ or ‘valued’ information can, in some circumstances, be somehow mapped to the binary one, thus

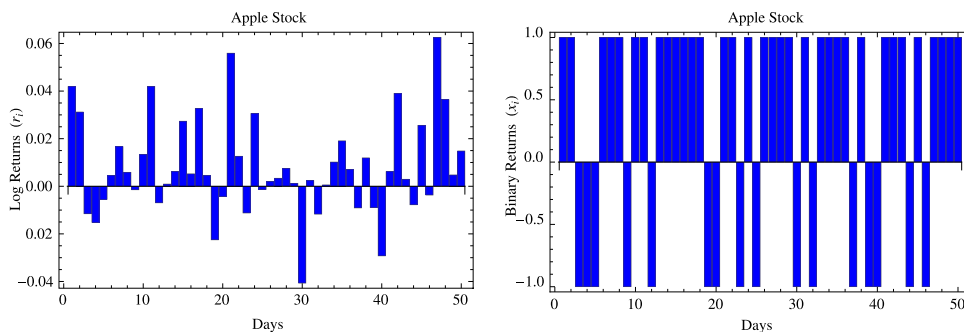
providing a robust, highly simplified, more easily modeled, and informative representation of the system.

Motivated by the above considerations, in this paper we further study, both empirically and theoretically, the relationship between weighted time series and their binary projections. We first provide robust empirical evidence of novel relationships between binary and non-binary properties of real financial time series. To this end, we use the daily closing prices of all stocks of three markets (S&P500, FTSE100 and NIKKEI225) over the period 2001–2011. We show that the average daily increment and average daily coupling of an empirical set of multiple time series are strongly and nonlinearly related to the corresponding average increment of the binary projections of the same time series. These empirical relations quantify in a novel way the strong correlations existing between the increments of individual stocks and the overall level of synchronization among all stocks in the market.

Building on this evidence, we then introduce a formalism to analytically characterize random ensembles of single and multiple time series with desired constraints. Specifically, we follow Jaynes' interpretation and re-derivation of statistical physics as an inference problem from partial macroscopic information to the unobservable microscopic configuration [18, 19]. We define statistical ensembles of matrices that maximize Shannon's entropy [20], subject to a set of desired constraints. This maximum-entropy approach is widely used in many areas, from neuroscience [21] to social network analysis [22] (and more recently network science in general [23]), where it is known under the name of exponential random graph (ERG) formalism. In the case of interest here, we introduce ensembles of maximum-entropy binary matrices that represent projections of single and multiple binary time series, subject to a set of desired constraints defined as simple empirical measurements. We discuss the main differences between our matrix ensembles and other techniques in time series analysis, including other ensembles of random matrices encountered in random matrix theory [24–28].

Our approach leads to a family of analytically solved null models that allow us to quantify the amount of information encoded in the chosen constraints, i.e. the selected observed properties of the binary projections of real time series. Different choices of the constraints lead to different stochastic processes, a result that allows us to relate known stochastic processes to the corresponding 'target' empirical properties defining the ensemble of time series spanned by the process itself. After applying the approach to the financial time series in our analysis, we compare the informativeness of various measured properties and show that different properties are more relevant for different time series and temporal windows. We also identify distinct regimes in the behaviour of multiple stocks and give the most likely explanation (endogenous, exogenous, or mixed) for the observed level of coordination or 'market mode', given the measured binary return at a given point in time. Finally, and most importantly, we show that our approach is able to reproduce and mathematically characterize the observed nonlinear relationships between binary and non-binary properties of real time series.

The rest of the paper is organized as follows. In section 2 we describe the data and provide empirical evidence of the relationships that motivate our work. In section 3 we introduce our theoretical formalism in its general form. In section 4 we apply the formalism to single time series, while in section 5 we apply it to single cross-sections (temporal snapshots) of multiple time series. Finally, in section 6 we consider our method in its full extent and apply it to entire spans of multiple time series, for different financial markets around the globe. We end with our conclusions in section 7.



**Figure 1.** ‘Weighted’ (left) versus ‘binary’ (right) time series of log-returns of the Apple stock over a period of 50 days starting from 7 May 2011.

## 2. Empirical results

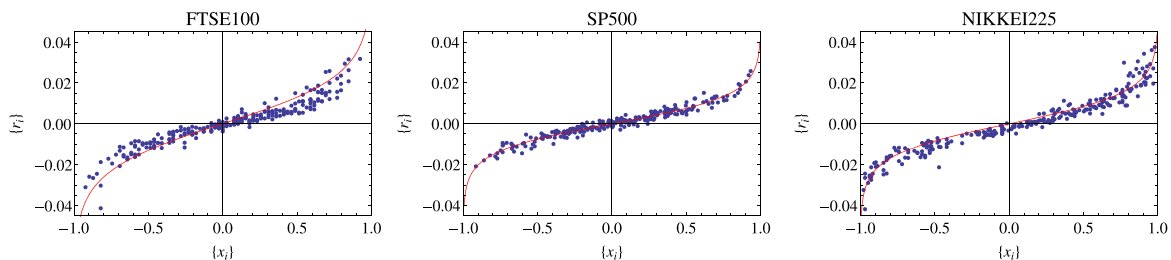
### 2.1. Data

We use daily closing prices, for the 10 year period ranging from 24 October 2001 to 18 October 2011, of all stocks from the indices S&P500, FTSE100 and NIKKEI225. For each index, we restrict our sample to the maximal group of stocks that are traded continuously throughout the selected period. This results in 445 stocks for the S&P500, 78 stocks for the FTSE100 and 193 stocks for the NIKKEI225.

We take logarithms of daily closing prices to obtain time series of *log-prices* that represent our original ‘signal’  $s_i(t)$ , where  $i$  labels stocks and  $t$  labels days in the sample. Correspondingly, we construct time series of *log-returns* where each entry represents the increment  $r_i(t)$  as defined in equation (1). Finally, we take the sign  $x_i(t)$  of each log-return  $r_i(t)$  to obtain an additional, binarized set of time series as in equation (2). We will refer to the binarized time series as the *binary projection* of the original time series. In figure 1 we show a simple example of a weighted time series, along with the corresponding binary projection. The (multiple) time series of  $r_i(t)$  and  $x_i(t)$  are the main objects of our analysis throughout the paper. Note that, while the use of log-returns rather than simple returns (i.e. price differences) in finance is an important step that allows the removal of overall trend effects over long time spans [5], the binary signature is actually independent of whether the original prices have been logarithmically transformed.

The main reason for choosing the daily frequency is to achieve an optimal level of structural compatibility between the data and the models we introduce later. As we discuss in detail in section 3, our models are binary, i.e. they only allow the two values  $\pm 1$  depending on whether the increment of the original time series is positive or negative. An increment of 0 is not admitted in the models: consistently, we chose a frequency for which zero increments are extremely rare in the data. In financial markets, this is the case for daily (or lower) frequency. Indeed, a zero return value occurs in less than 0.2% of the cases in our daily data (when this happens, we randomly switch the corresponding binary increment to either +1 or  $-1$  with equal probability). Higher-frequency data feature an increasing percentage of zero returns, a property that calls for an extension of the models considered here.

It should be noted that other types of binary time series, different from the  $\pm 1$  type considered here, can also be defined. Most notably, 0/1 binary time series can indicate the occurrence of an event in a time period, i.e. whether the event happened (1) or not (0). Financial



**Figure 2.** Nonlinear relationship between the average daily increment (weighted return) and the average daily sign (binary return) over all stocks in the FTSE100 (left), S&P500 (center) and NIKKEI225 (right) in various years (2003, 2007, and 2004, respectively). Here each point corresponds to one day in the time interval of 250 trading days (approximately one year). The red line represents the best fit with the function  $y = a \cdot \operatorname{artanh}x$ , whose use is theoretically justified later in section 6.

examples include time series of recession indicators [29, 30] or of ‘switching points’ in stock returns [31]. For such 0/1 binary time series, correlations may not be very informative when measuring a dependence between the dichotomous variables. To confront this gap, in recent years new methods have been introduced, like the auto-persistence function and auto-persistence graph [29]. In these methods, the dependence structure among the observations is described in terms of conditional probabilities, rather than correlations. Although throughout this paper we will be entirely focusing on  $\pm 1$  binary time series that naturally descend from the original *signed* time series of fluctuations, it is interesting to notice that our approach can be extended, with slight modifications, to 0/1 time series as well. To this end, one needs to re-express all quantities in terms of a 0/1 binary variable  $y \equiv (x + 1)/2$ , where  $x$  is our  $\pm 1$  binary variable, and adapt our approach accordingly.

## 2.2. Nonlinear binary/non-binary relationships

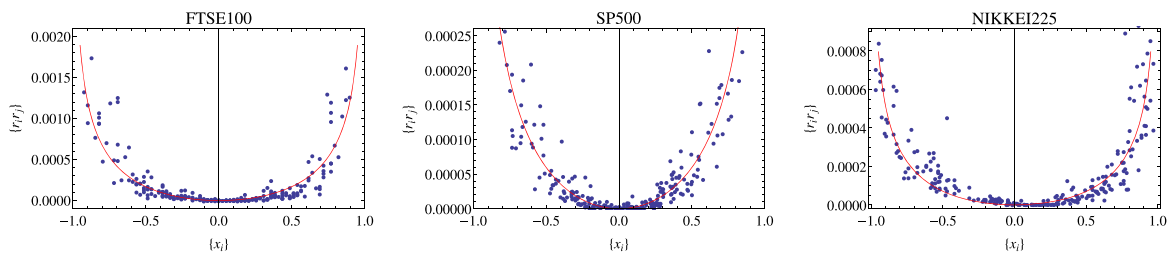
We now come to the main empirical findings that motivate our paper. For each index and for each day  $t$  in the sample, we first calculate the average (over all stocks) weighted return, that we denote as  $\{r_i(t)\}$  and define as

$$\{r_i(t)\} \equiv \frac{1}{N} \sum_{i=1}^N r_i(t). \quad (3)$$

Note that the above expression does not depend on the particular stock  $i$ , but it does depend on time  $t$ . Our unconventional choice of the symbol  $\{\cdot\}$  to denote an average over stocks is to avoid confusion with temporal averages, which will be denoted by the more usual bar ( $\bar{\cdot}$ ) later in the paper. Similarly, we calculate the corresponding average binary return  $\{x_i(t)\}$ , defined as

$$\{x_i(t)\} \equiv \frac{1}{N} \sum_{i=1}^N x_i(t). \quad (4)$$

In figure 2 we plot  $\{r_i(t)\}$  as a function of  $\{x_i(t)\}$  for all days of various 1 year intervals and for the three indices separately. We find a strong nonlinear dependency between the two quantities. Note that the average binary return is bound between  $-1$  and  $+1$  by construction, but the average weighted return is unbounded from both sides. While there are in principle infinite



**Figure 3.** Nonlinear relationship between the average daily coupling (weighted coupling) and the average daily sign (binary return) over all stocks in the FTSE100 (left), S&P500 (center) and NIKKEI225 (right) in various years (2003, 2007, and 2004, respectively). Here each point corresponds to one day in the time interval of 250 trading days (approximately one year). The red line represents the best fit with the function  $y = b \cdot (\text{artanh}x)^2$ , whose use is theoretically justified later in section 6.

values of  $\{r_i(t)\}$  that are consistent with the same value of  $\{x_i(t)\}$ , we observe a tight relationships between the two quantities. This relationship can be fitted by a one-parameter curve of the form

$$\{r_i(t)\} = a \cdot \text{artanh}\left[\{x_i(t)\}\right] = \frac{a}{2} \ln \frac{1 + x_i(t)}{1 - x_i(t)} \quad (5)$$

(the theoretical justification for this functional form will be given in section 6), where  $a$  is in general different for different years and different indices. Still, as we show later, for a given year and market the average weighted return of any day  $t$  is to a large extent predictable (out of sample) from the average binary return of the same day, once  $a$  is known (for instance by fitting the above curve to the data for a past time window). In section 6 we will also show that the nonlinear character of the observed relations is a genuine signature of correlation in the data, as an uncorrelated null model shows a completely linear behaviour.

There is another empirical relationship, involving a higher-order quantity. For each index and for each day  $t$  in the sample, we calculated what we will call the average ‘coupling’ over the  $N(N - 1)/2$  distinct pairs of stocks:

$$\{r_i(t)r_j(t)\} \equiv \frac{2\sum_{i<j} r_i(t)r_j(t)}{N(N - 1)} \quad (6)$$

(so now the symbol  $\{\cdot\}$  indicates an average over *pairs* of stocks). In figure 3 we plot  $\{r_i(t)r_j(t)\}$  as a function of the average binary return  $\{x_i(t)\}$ , for the same data as in figure 2. Again, we find a strong nonlinear dependency, where for a given value of the average binary return of day  $t$  there is a typical value of the average coupling among all stocks in the same day. The relationship can be fitted by a one-parameter curve that diverges at  $\{x_i\} = \pm 1$ . As we show in section 6, an uncorrelated null model would yield a different, parabolic curve with no divergences. Again, this means that the empirical trend is due to genuine correlations, whose nature will be clarified later on in the paper.

There are even more examples of dependencies that we can find between binary and non-binary properties in the data. However, in one way or another all these relationships, including that shown in figure 3, ultimately derive from equation (5). For this reason, we refrain from showing redundant results and focus on the empirical findings discussed so far.



The above analysis indicates that the binary signature of financial time series contains relevant information about the original data. While the binary signature is *a priori* a many-to-one projection involving a significant information loss, we empirically find that there are properties (namely the average return and average coupling) for which the projection is virtually a one-to-one ‘quasi-stationary’ transformation (on appropriate time scales, as we show in section 6), allowing the reconstruction of the corresponding original, weighted properties to a great extent. Rather than exploring the practical aspects of this possibility of reconstruction of the original signal from its binary projection, in this paper we are interested in understanding the origin of this behaviour and providing a simple data-driven model of it. This will be ultimately achieved in section 6, where we also show that the binary/non-binary relations we have documented are a novel quantification of the fact that extreme price increments occur more often when most stocks move in the same direction. This is an important type of correlation between the magnitude of log-returns of individual time series and the level of synchronization (common sign) of the increments of all stocks in the market.

### 3. Maximum-entropy matrix (MEM) ensembles

Having established that the binary projections of real time series contain non-trivial information, in the rest of the paper we introduce a theory of binary time series aimed, among other things, at reproducing the observed nonlinear relationships showed in figures 2 and 3. In our approach, we regard a synchronous set of binary time series as a  $\pm 1$  matrix and we introduce an ensemble of such matrices via the maximization of Shannon’s entropy, subject to the constraint that some specified properties of the ensemble match their observed values. An analogous approach is widely used, e.g., in network analysis and known under the name ERG [23]. Moreover, we provide an analytical maximum-likelihood method to find the optimal values of the parameters governing the ensembles, which is again similar in spirit to a method that has been recently introduced for networks [32–34]. Finally, we describe Akaike’s information criterion (AIC) [35], which we will use to rank and compare the performance of different null models when fitted to the same data.

Being entropy-based, our approach automatically allows us to measure the amount of information encoded into the observed properties chosen as constraints, i.e. how much information is gained about the original (set of) time series once those properties are measured. It also allows us to identify, given a set of measured properties, which ones are more informative and which ones can be discarded, as we show on specific financial examples. Our framework turns out to reproduce the observed nonlinear relationships very well, thus providing a simple mathematical explanation and functional form for the plots shown in the previous section. Moreover, we are able to identify, as a function of the binary return only, distinct regimes in the collective behavior of stocks, namely a ‘coordinated’ regime dominated by market-wide interactions, an ‘uncoordinated’ regime dominated by stock-specific noise and an ‘intermediate’ regime where both market-wide and stock-specific information is relevant.

We incidentally note that, despite the available variety of refined and advanced techniques in time series analysis [36], how one can quantify (in the sense of statistical ensembles) how much information is actually encoded into any given, measurable property of a time series is still not fully understood. While most studies, starting from the celebrated work by Kolmogorov about the algorithmic complexity of sequences of symbols [37], have addressed the



quantification of the information content of a single time series, much less is known about the information encoded in the measured value of a given time series property (which, necessarily, involves the idea of an entire *ensemble* of time series consistent with the measured value itself). Our approach can provide an answer to such a question, by associating an absolute level of uncertainty (entropy) to each observable of an empirical (set of) time series. In relative terms, this also allows us to compare the information content of different properties of a time series, thereby indicating which measured property is the most informative about the original time series.

As a final consideration, it is worth mentioning that the MEM ensembles that we introduce are clearly related to (and, depending on the specification, potentially overlapping with) some ensembles that are well studied by random matrix theory [38–43]. However, our approach is different since we generate ensembles of matrices whose probability distributions are determined by the kind of partial information (empirically measured constraint) about the real system. In this approach the maximization of Shannon’s entropy, given some real-world available information, yields the least biased probability distribution (over the space of possible matrices) consistent with the data. This formalism allows us to relate the probabilistic structure of each matrix ensemble with the choice of the original observed property, or constraint. Similarly, since our matrices represent (multiple) time series, we are able to connect the various ensembles to simple stochastic processes induced by the associated matrix probabilities and, again, to the chosen empirical property specifying the ensembles themselves.

### 3.1. Exponential random matrices (ERMs)

We first analytically characterize the properties of families of randomized matrices. More generally, we introduce a matrix ensemble that maximizes Shannon’s entropy, while enforcing a set of observed constraints (selected time series properties). This procedure is analogous to e.g., that leading to the definition of ERGs in network theory [23]. However, we will modify it to accommodate  $\pm 1$  matrices, as opposed to 0/1 or non-negative matrices that describe binary and weighted networks, respectively. The resulting ensemble can thus be denoted as the MEM ensemble or equivalently the ERMs model.

Let us consider the ensemble of all  $\pm 1$  matrices with dimensions  $N \times T$ . Each such matrix can represent  $N$  synchronous time series, all of duration  $T$  (for instance, if applied to a set of multiple financial time series,  $N$  refers to the number of stocks and  $T$  to the number of time steps). Let  $\mathbf{X}$  denote a generic matrix in the ensemble, and  $x_i(t)$  its entry ( $1 \leq i \leq N, 1 \leq t \leq T$ ). Let  $\mathbf{X}^*$  be the particular real matrix that we observe. In other words, our ensemble is composed of all possible matrices  $\mathbf{X}$  of the same type as  $\mathbf{X}^*$ , and includes  $\mathbf{X}^*$  itself. For any data-dependent property  $R$ , we will consider the value  $R(\mathbf{X})$  obtained when  $R$  is measured on the particular matrix  $\mathbf{X}$ . For each matrix  $\mathbf{X}$  in the ensemble, we will assign an occurrence probability  $P(\mathbf{X})$ . The expectation value (ensemble average) of a property  $R$  can be expressed as

$$\langle R \rangle = \sum_{\mathbf{X}} R(\mathbf{X}) P(\mathbf{X}), \quad (7)$$

where the sum runs over all matrices in the ensemble.

At this point, we introduce a set of constraints denoted by the vector  $\vec{C}$ . The constraints are meant to ensure that a given set of observed properties  $\vec{C}(\mathbf{X}^*)$  in the real matrix  $\mathbf{X}^*$  is reproduced by the ensemble itself. In our method we will enforce ‘soft’ constraints by requiring

that their expectation value  $\langle \vec{C} \rangle$  equals the observed one. The resulting ensemble is a *canonical* one where each matrix  $\mathbf{X}$  is assigned a probability  $P(\mathbf{X})$  that maximizes Shannon's entropy

$$S \equiv - \sum_{\mathbf{X}} P(\mathbf{X}) \ln P(\mathbf{X}) \quad (8)$$

subject to the normalization constraint

$$\sum_{\mathbf{X}} P(\mathbf{X}) = 1 \quad (9)$$

and to the chosen vector of constraints

$$\langle \vec{C} \rangle = \sum_{\mathbf{X}} C(\mathbf{X}) P(\mathbf{X}) = \vec{C} \quad (10)$$

that we enforce in order to reproduce the desired set of observed quantities.

The solution to the above constrained maximization problem is standard (see for instance [23] for a recent derivation in the context of networks). We first introduce the Lagrange multipliers  $\alpha$  and  $\vec{\theta}$ , enforcing equations (9) and (10) respectively, and then require that the functional derivative of Shannon's entropy (plus the constraining terms) vanishes:

$$\frac{\partial}{\partial P(\mathbf{X})} \left\{ S + \alpha \left[ 1 - \sum_{\mathbf{X}} P(\mathbf{X}) \right] + \sum_i \theta_i \left[ C_i - \sum_{\mathbf{X}} C(\mathbf{X}) P(\mathbf{X}) \right] \right\} = 0.$$

This yields

$$\ln P(\mathbf{X}) + 1 + \alpha + \sum_i \theta_i C_i(\mathbf{X}) = 0 \quad (11)$$

for any matrix  $\mathbf{X}$ . Using a notation that makes the dependence of all quantities on  $\vec{\theta}$  explicit, we then obtain

$$P(\mathbf{X} | \vec{\theta}) = \frac{e^{-H(\mathbf{X}, \vec{\theta})}}{Z(\vec{\theta})}, \quad (12)$$

where  $H(\mathbf{X}, \vec{\theta})$  is the *Hamiltonian*

$$H(\mathbf{X}, \vec{\theta}) \equiv \vec{\theta} \cdot \vec{C}(\mathbf{X}) = \sum_i \theta_i C_i(\mathbf{X}), \quad (13)$$

which is a linear combination of the constraints, and  $Z(\vec{\theta})$  is the *partition function*

$$Z(\vec{\theta}) \equiv e^{\alpha+1} = \sum_{\mathbf{X}} e^{-H(\mathbf{X}, \vec{\theta})}, \quad (14)$$

which is the normalizing constant for the probability. Consistently, we can rewrite equation (7) more explicitly as a function of  $\vec{\theta}$ :

$$\langle R \rangle_{\vec{\theta}} \equiv \sum_{\mathbf{X}} R(\mathbf{X}) P(\mathbf{X} | \vec{\theta}), \quad (15)$$

where  $\langle \cdot \rangle_{\vec{\theta}}$  indicates that the ensemble average is evaluated at the particular parameter value  $\vec{\theta}$ .

Equations (12)–(14) define the MEM or ERM model. Specifically, the model yields the probability distribution over a specified ensemble of matrices, which maximizes the entropy under a set of generic constraints. The guiding principle is that the probability distribution (over

microscopic states) which have maximum entropy, subject to observed (macroscopic) properties, provides the most unbiased representation of our knowledge of the state of a system [19]. To put it in a more physical frame, this is analogous to the Gibbs–Boltzmann distribution over the microstates of a large system at a well defined temperature, given the thermodynamic (macroscopic) observables such as the total energy.

### 3.2. Maximum-likelihood parameter estimation

The above derivation shows that the expectation value of any property of the ensemble depends *functionally* on the specific enforced constraints  $\vec{C}$  through the resulting structure of  $P(\mathbf{X} | \vec{\theta})$ . Of course, it also depends *numerically* on the measured values  $\vec{C}(\mathbf{X}^*)$  of the constraints themselves, through the particular parameter value (that we denote by  $\vec{\theta}^*$ ) required in order to enforce that the expected and observed values of  $\vec{C}$  match:

$$\langle \vec{C} \rangle_{\vec{\theta}^*} = \vec{C}(\mathbf{X}^*). \quad (16)$$

We now show that the value  $\vec{\theta}^*$  that satisfies equation (16) coincides with the value that maximizes the likelihood to generate the empirical data, as in the corresponding maximum likelihood (ML) approach to network ensembles [32, 44].

We start by writing the log-likelihood function of an observed matrix  $\mathbf{X}^*$  generated by the parameters  $\vec{\theta}$ :

$$\lambda(\vec{\theta}) \equiv \ln P(\mathbf{X}^* | \vec{\theta}) = -H(\mathbf{X}^*, \vec{\theta}) - \ln Z(\vec{\theta}). \quad (17)$$

We then look for the particular value  $\vec{\theta}^*$  that maximizes  $\lambda(\vec{\theta})$ , i.e.

$$\vec{\nabla} \lambda(\vec{\theta}^*) = \left[ \frac{\partial \lambda(\vec{\theta})}{\partial \vec{\theta}} \right]_{\vec{\theta}=\vec{\theta}^*} = \vec{0} \quad (18)$$

(it is easy to check that the higher-order derivative confirms that  $\vec{\theta}^*$  is a point of maximum). This leads to

$$\vec{\nabla} \lambda(\vec{\theta}^*) = \left[ -\vec{C}(\mathbf{X}^*) - \frac{1}{Z(\vec{\theta})} \frac{\partial Z(\vec{\theta})}{\partial \vec{\theta}} \right]_{\vec{\theta}=\vec{\theta}^*} = \vec{0} \quad (19)$$

the solution for that yields the ML condition

$$\vec{C}(\mathbf{X}^*) = \sum_{\mathbf{X}} \frac{\vec{C}(\mathbf{X}) e^{-H(\mathbf{X}, \vec{\theta}^*)}}{Z(\vec{\theta}^*)} = \langle \vec{C} \rangle_{\vec{\theta}^*}, \quad (20)$$

which coincides with equation (16). Thus the likelihood of the real matrix  $\mathbf{X}^*$  is maximized by the specific parameter choice such that the ensemble average of each constraint equals its empirical value measured on  $\mathbf{X}^*$ , automatically ensuring that the desired constraints are met.

### 3.3. Model selection

We finally show how we can use AIC to rank the performance of different models, i.e. different choices of the constraints, in reproducing the same data. The AIC is an information-theoretic measure of the relative goodness of fit of a model, as compared to a set of alternative models all used to explain the same data [35]. It offers a relative measure of the information lost when the given model is used to describe reality. The power of AIC (and other similar criteria [45]) lies in the possibility to rank a set of models in terms of their achieved trade-off between accuracy (good fit to the data) and parsimony (low number of free parameters) [45]. In general, for the  $k$ th model in a set of selected models, AIC is defined as

$$AIC_k = 2n_k - 2\lambda_k^* \quad (21)$$

where  $n_k$  is the number of free parameters in the  $k$ th model and  $\lambda_k^*$  is the maximized log-likelihood of the data under the same model. The above expression effectively discounts the number  $n_k$  of parameters (complexity) from the maximized likelihood  $\lambda_k^*$  (accuracy). The model with the lowest value of  $AIC_k$  (let us denote this value by  $AIC_{\min}$ ) is the ‘best’ model in the considered set, achieving the optimal trade-off [45].

In the ERM/MEM family of models we have introduced, a model is uniquely specified by the choice of the constraints  $\vec{C}$ . Given a  $N \times T$  data matrix  $\mathbf{X}^*$  and a set  $\{\vec{C}_1, \dots, \vec{C}_m\}$  of  $m$  possible choices of constraints, each of the resulting  $m$  models has an AIC value

$$AIC_k = 2n_k - 2 \ln P_k(\mathbf{X}^* | \vec{\theta}_k^*) \quad k = 1, m \quad (22)$$

where  $n_k$  is the dimensionality of the vector  $\vec{C}_k$ ,  $\ln P_k(\mathbf{X}^* | \vec{\theta}_k^*)$  is the maximized log-likelihood of model  $k$ , and  $\vec{\theta}_k^*$  is the parameter value maximizing such log-likelihood. Within our framework, AIC identifies which measured property  $\vec{C}_k(\mathbf{X}^*)$  is most informative about the entire time series  $\mathbf{X}^*$ .

In order to understand whether models with values of AIC larger than but close to  $AIC_{\min}$  are still competitive, it is customary to define the so-called ‘AIC weights’ which provide a normalized strength of evidence for a model [45]. For each model  $k$  in the set of  $m$  models, one first calculates the difference  $\Delta_k = AIC_k - AIC_{\min}$  and then defines the AIC weight

$$w_k \equiv \frac{e^{-\Delta_k/2}}{\sum_{r=1}^m e^{-\Delta_r/2}}. \quad (23)$$

The AIC weight  $w_k$  represents the probability that the  $k$ th model is the best one among the  $m$  selected models. For instance, an AIC weight of  $w_k = 0.75$  indicates that, given the data, model  $k$  has a 75% chance of being the best model among the  $m$  candidate ones. If two or more models have comparable AIC weights (e.g.  $w_1 = 0.6$ ,  $w_2 = 0.4$  or  $w_1 = 0.35$ ,  $w_2 = 0.25$ ,  $w_3 = 0.4$ ), then there is no evidence that the model with the highest AIC weight (lowest AIC value) is clearly outperforming the other ones. All the models with comparable weights should be considered as competing alternatives, in principle leading to the problem of multi-model inference [45].

## 4. Single time series

In this section we consider the first family of specifications of our general approach outlined in section 3. We focus on the simple case of single time series ( $N = 1$ ), where the ensemble of

$N \times T$  matrices reduces to an ensemble of  $1 \times T$  matrices, or equivalently of  $T$ -dimensional row vectors. Each such vector will still be denoted by  $\mathbf{X}$ . We assume long time series, i.e.  $T \gg 1$ .

This first specification of our abstract formalism is not meant to provide realistic models for the evolution of the binary increments of real financial time series. Rather, it allows us to make different sorts of considerations. On one hand, it allows us to introduce our formalism using simpler examples first, establishing the basis for the more general cases (leading to the main results of the paper) that will be introduced later. On the other hand, it emphasizes that different and well known (one-dimensional) stochastic processes are found as particular examples of maximum-entropy ensembles defined by specific constraints that are otherwise obscure. Identifying these ‘driving constraints’ underlying common stochastic processes will help us interpret such processes in the light of the empirical properties being reproduced. Finally, our approach allows us to identify, given the data and given a set of simple properties, which of these properties is encoding the largest amount of information about the original binary signature.

Let  $\mathbf{X}$  denote a single time series with entries  $x(t)$ , where  $1 \leq t \leq T$ , each representing a temporal increment. We will denote the average increment (first moment) as

$$M_1(\mathbf{X}) \equiv \overline{x(t)} = \frac{1}{T} \sum_{t=1}^T x(t). \quad (24)$$

Note that the second moment is always

$$M_2(\mathbf{X}) \equiv \overline{x^2(t)} = \frac{1}{T} \sum_{t=1}^T x^2(t) = 1, \quad (25)$$

so the sample variance is

$$M_2(\mathbf{X}) - M_1^2(\mathbf{X}) = 1 - \overline{x(t)}^2. \quad (26)$$

We also define the  $\tau$ -delayed product (with  $0 \leq \tau \leq T$ )

$$B_\tau(\mathbf{X}) \equiv \overline{x(t) \cdot x(t + \tau)} = \frac{1}{T} \sum_{t=1}^T x(t) \cdot x(t + \tau) \quad (27)$$

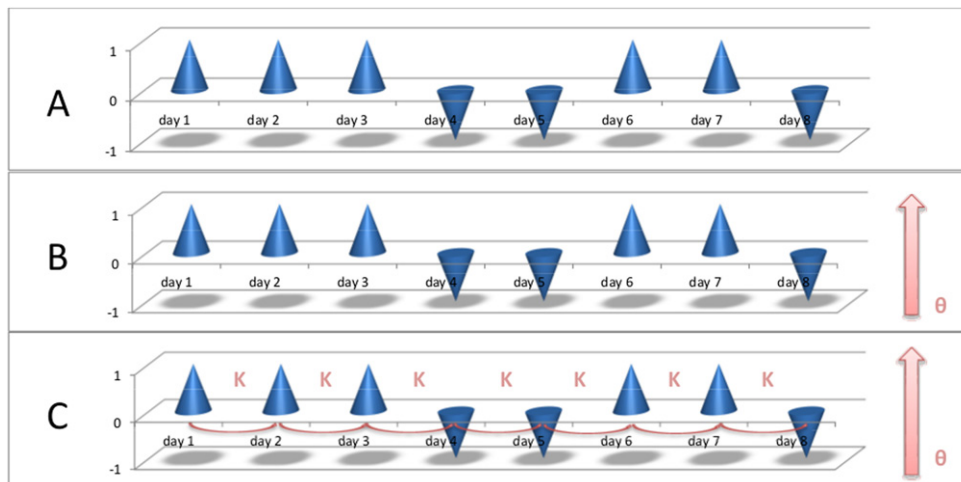
where we have introduced periodic boundary conditions:

$$x(T + \tau) \equiv x(\tau) \quad \text{with} \quad 0 \leq \tau \leq T. \quad (28)$$

The above periodicity condition is inessential, since we could have used a definition avoiding its introduction, but it makes some expressions simpler in what follows. Periodicity implies that the normalized (between  $-1$  and  $+1$ ) autocorrelation function (with delay  $\tau$ ) can be defined as

$$\begin{aligned} A_\tau(\mathbf{X}) &\equiv \frac{\overline{x(t) \cdot x(t + \tau)} - \overline{x(t)} \cdot \overline{x(t + \tau)}}{\overline{x^2(t)} - \overline{x(t)}^2} \\ &= \frac{B_\tau(\mathbf{X}) - M_1^2(\mathbf{X})}{1 - M_1^2(\mathbf{X})}. \end{aligned} \quad (29)$$

Since a  $(\pm 1)$  binary time series can also be regarded as a chain of classical spins pointing either up or down, it is straightforward to consider simple, analytically solved spin models as the starting point, since these models are defined in terms of a ‘physical’ Hamiltonian that has



**Figure 4.** Illustration of our mapping from single binary time series to spin models. Each time series is regarded as a chain of  $\pm 1$  spins, where the value of the spin indicates if the daily return of the stock is positive (+1) or negative (-1). In each model we enforce different constraints that imply different spin models and different stochastic processes. Given the same time series, we consider three possible models. (A) We enforce no constraint, which translates into a chain of non-interacting spins without external field (uniform random walk). (B) We enforce the total temporal increment, which translates into a chain of non-interacting spins with external field (biased random walk). (C) We enforce both the total increment and the one-lagged autocorrelation, which translates into a chain of spins with first-neighbour interactions and external field (Markov process).

precisely the same structure of our ‘information-theoretic’ Hamiltonian defined in equation (13). In what follows, we introduce various model specifications. For each model, we introduce the constraints that we enforce and the resulting Hamiltonian as described in section 3.1. Different constraints correspond to different spin models and lead to different stochastic processes. This is pictorially illustrated in figure 4. The free parameters conjugated to the constraints will be fitted according to the ML principle described in section 3.2. Different models will be ranked according to the AIC weights introduced in section 3.3.

#### 4.1. Uniform random walk

The most trivial model is one where we enforce no constraint, i.e. there is no free parameter and the Hamiltonian is

$$H(\mathbf{X}) = 0. \quad (30)$$

Physically, the above Hamiltonian describes a gas of  $T$  non-interacting ‘spins’ in a vacuum, i.e. in absence of an external magnetic field. This model is discussed in the appendix. The probability of the occurrence of a time series  $\mathbf{X}$  is completely uniform over the ensemble of all binary time series of length  $T$ . All the  $T$  elements of  $\mathbf{X}$  are mutually independent and identically distributed. This results in a completely uniform random walk with zero expected value for each increment:

$$\langle x(t) \rangle = 0. \quad (31)$$

While the (ensemble) variance of each increment equals

$$\text{Var}[x(t)] \equiv \langle x^2(t) \rangle - \langle x(t) \rangle^2 = 1. \quad (32)$$

This trivial model generates a symmetric random walk. Since the expected return is zero, and the uncertainty is maximal, the variance is also maximal (for a  $\pm 1$  binary random variable). Financially, the model assumes that the stock fluctuates randomly, with no memory, and with no overall ‘price drift’. This is the most basic model of price dynamics that has been considered in the financial literature since the pioneering work of Bachelier [1], here adapted to the case of binary time series.

The model can be used as a basic benchmark for checking the performance of our other models. This comparison will be studied in section 4.4. Since here the likelihood is independent of any parameter, the AIC of the model can be calculated using equation (22) where the probability is given by equation (A.3) (see appendix) and the number of parameters is  $n_k = 0$ .

#### 4.2. Biased random walk

We now consider the total increment as the simplest non-trivial (one-dimensional) constraint:

$$C(\mathbf{X}) = T \cdot M_1(\mathbf{X}) = T \cdot \overline{x(t)}. \quad (33)$$

This leads to the Hamiltonian

$$H(\mathbf{X}, \theta) = \theta \sum_{t=1}^T x(t), \quad (34)$$

which coincides with the physical Hamiltonian for a gas of  $T$  non-interacting ‘spins’ in a common external ‘magnetic field’  $-\theta$ .

As we show in the appendix, this model generates a *biased* random walk where the probability  $P_t(x | \theta)$  of a given increment  $x = \pm 1$  at time  $t$  is

$$P_t(x | \theta) = \frac{e^{-\theta x}}{e^{-\theta} + e^{+\theta}}. \quad (35)$$

The expected return is the hyperbolic tangent

$$\langle x(t) \rangle_\theta = -\tanh \theta, \quad (36)$$

while the variance is

$$\text{Var}[x(t)] = 1 - \tanh^2 \theta. \quad (37)$$

Financially, this model still assumes no memory in the fluctuations of a given stock, but it introduces a ‘price drift’ in terms of a non-zero expected return.

The maximum likelihood condition (16), fixing the value  $\theta^*$  of the parameter  $\theta$  given a real time series  $\mathbf{X}^*$ , leads to

$$\theta^* = -\frac{1}{2} \ln \left[ \frac{1 + \overline{x^*(t)}}{1 - \overline{x^*(t)}} \right]. \quad (38)$$



The maximized likelihood for the model is

$$P(\mathbf{X}^* | \theta^*) = \prod_{t=1}^T P_t(x^*(t) | \theta^*) \quad (39)$$

which, using equation (22) with  $n_k = 1$ , can be used to measure the AIC (see section 3.3) of the model, based on the observed data. This will be done in section 4.4.

### 4.3. One-lagged model

Let us now explore a more complex model of collective behavior. The models considered so far were non-interacting, i.e. each return in the time series was independent of the previous outcomes. Now we consider a model where, besides the constraint on the total increment specified in equation (33), we enforce an additional constraint on the time-delayed (lagged) quantity  $T \cdot B_1(\mathbf{X})$ , where  $B_1(\mathbf{X})$  is defined in equation (27) with  $\tau = 1$ . Financially, this amounts to enforcing the average return *and* the average one-step temporal autocorrelation of the time series. In other words, besides a price drift, we also introduce a short-term memory.

The resulting two-dimensional constraint can be written as

$$\vec{C}(\mathbf{X}) = \begin{pmatrix} C_1(\mathbf{X}) \\ C_2(\mathbf{X}) \end{pmatrix} = T \cdot \begin{pmatrix} M_1(\mathbf{X}) \\ B_1(\mathbf{X}) \end{pmatrix}. \quad (40)$$

If we write the corresponding Lagrange multiplier as

$$\vec{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = - \begin{pmatrix} I \\ K \end{pmatrix}, \quad (41)$$

then the Hamiltonian reads

$$H(\mathbf{X}, I, K) = -I \sum_{t=1}^T x(t) - K \sum_{t=1}^T x(t)x(t+1), \quad (42)$$

where we consider a periodicity condition as in equation (28) with  $\tau = 1$ , i.e.  $x(T+1) \equiv x(1)$ . Note that, when  $\mathbf{X}$  is a real binary time series of length  $T$ , this condition can be always enforced by adding one last (fictitious) timestep  $T+1$  and a corresponding increment  $x(T+1)$  chosen equal to  $x(1)$ . For long time series (large  $T$ ), the effects induced by this addition are negligible.

The above Hamiltonian coincides with that for the one-dimensional Ising model with periodic boundary conditions [46], which is a model of interacting spins under the influence of an external ‘magnetic’ field  $I$ . The model is analytically solvable (see the appendix for the complete derivation), which allows us to apply it to real time series in our formalism. In our setting, each time step  $t$  is seen as a site in an ordered chain of length  $T$ , and each value  $x(t) = \pm 1$  is seen as the value of a spin sitting at that site. ‘First-neighbour interactions’ along the chain of spins are here interpreted as one-lagged memory effects. As a result of these interactions, the model generates time series according to a Markov process where the probability of an increment  $x(t+1)$  depends on the realized increment  $x(t)$  at the previous time step  $t$ . This is evident from the solution of the model, see e.g. equation (A.32) in the appendix.

The solution of the model yields the following expectation values

$$\langle M_1 \rangle_{I,K} = \frac{e^{2K} \sinh I}{\sqrt{1 + e^{4K} \sinh^2 I}} \quad (43)$$

$$\langle B_\tau \rangle_{I,K} = \frac{e^{4K} \sinh^2 I + (\lambda_1/\lambda_2)^\tau}{1 + e^{4K} \sinh^2 I} \quad (44)$$

(see the appendix) where

$$\lambda_1 = e^K \cosh I + \sqrt{e^{2K} \sinh^2 I + e^{-2K}}, \quad (45)$$

$$\lambda_2 = e^K \cosh I - \sqrt{e^{2K} \sinh^2 I + e^{-2K}}. \quad (46)$$

The resulting expected value of the normalized autocorrelation defined in equation (47) is simply

$$\langle A_\tau \rangle_{I,K} = \left( \frac{\lambda_1}{\lambda_2} \right)^\tau. \quad (47)$$

The above expressions allow us to calculate all the relevant expected properties of the time series generated by the model, once the parameters  $I$  and  $K$  are set to the values  $I^*$  and  $K^*$  maximizing the likelihood  $P(\mathbf{X}^* | I, K)$  of the observed time series  $\mathbf{X}^*$ . These values are the solutions of the coupled equations

$$M_1(\mathbf{X}^*) = \langle M_1 \rangle_{I^*,K^*} \quad (48)$$

$$B_1(\mathbf{X}^*) = \langle B_1 \rangle_{I^*,K^*}, \quad (49)$$

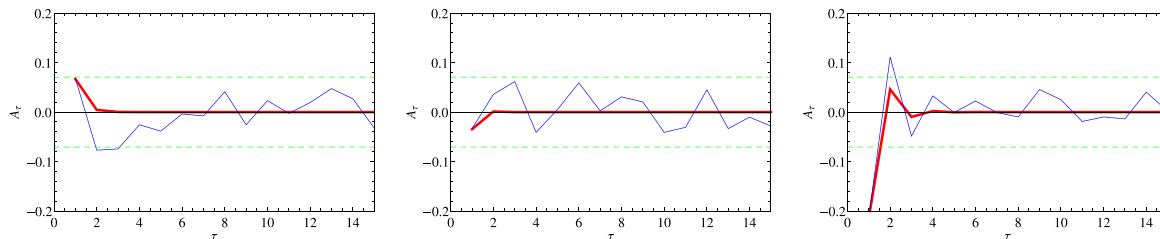
where  $M_1(\mathbf{X}^*)$  and  $B_1(\mathbf{X}^*)$  are the empirical values measured on the real data  $\mathbf{X}^*$ . The maximized likelihood of the model can be calculated as  $P(\mathbf{X}^* | I^*, K^*)$ , where  $P(\mathbf{X} | I, K)$  is given by equation (A.32) in the appendix. From the maximized likelihood, the AIC can be easily obtained using equation (22) with  $n_k = 2$ .

Note that the values  $I^*$  and  $K^*$  are such that the first point of the expected autocorrelation function,  $\langle A_1 \rangle_{I^*,K^*}$ , is necessarily equal to the observed value  $A_1(\mathbf{X}^*)$ . Based on this first value alone, the model will provide the full expected autocorrelation  $\langle A_\tau \rangle_{I^*,K^*}$  as follows:

$$\langle A_\tau \rangle_{I^*,K^*} = \left( \frac{\lambda_1}{\lambda_2} \right)_{I^*,K^*}^\tau = \left[ A_1(\mathbf{X}^*) \right]^\tau. \quad (50)$$

Comparing the above expression, for  $\tau > 1$ , with the observed autocorrelation function  $A_\tau(\mathbf{X}^*)$  is an important test of the model. Note that, since  $-1 \leq A_1(\mathbf{X}^*) \leq +1$ , the absolute value of the autocorrelation function  $\langle A_\tau \rangle_{I^*,K^*}$  is necessarily decreasing. If  $A_1(\mathbf{X}^*) > 0$  then  $\langle A_\tau \rangle_{I^*,K^*}$  will be positive (and exponentially decreasing) for all values of  $\tau$ . By contrast, if  $A_1(\mathbf{X}^*) < 0$  then  $\langle A_\tau \rangle_{I^*,K^*}$  will be an oscillating function (modulated by a decreasing exponential), and will take negative values when  $\tau$  is odd and positive values when  $\tau$  is even.

In figure 5 we compare the measured autocorrelation, equation (29), with the predicted one, equation (50), for three different S&P500 stocks (USB, Qcom, and MJN) over a period of



**Figure 5.** Measured autocorrelation (blue) of three different S&P500 stocks (Qcom, USB, and MJN respectively) over a period of 800 trading days (approximately 3.5 years), and comparison with the predicted autocorrelation  $\langle A_\tau \rangle_{I,K}$  generated by the one-lagged (one-dimensional Ising) model (red). The green lines represent the noise level, calculated as  $\pm 2$  standard deviations of the Fisher-transformed autocorrelation.

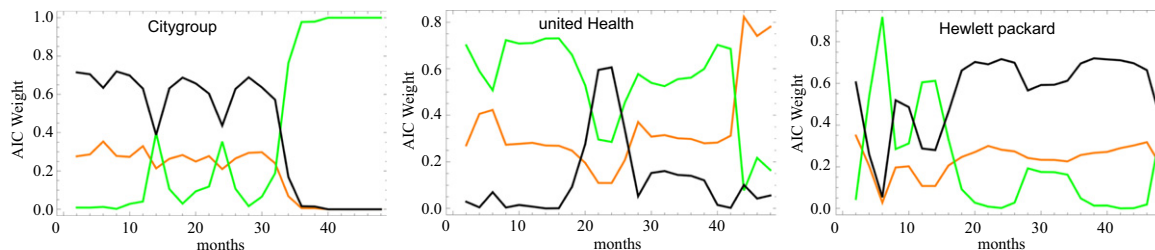
800 trading days (approximately 3.5 years). As expected, we see that the first point (one-lagged autocorrelation) is always reproduced exactly. We also confirm that, depending on the sign of the first point, the predicted trend is either exponentially decreasing (e.g. for the USB stock on the left) or oscillating (e.g. the Qcom and MJN stocks). The dashed lines indicate the noise level, which we arbitrarily fixed at two standard deviations of the fisher-transformed<sup>1</sup> autocorrelation. The behaviour of the USB and Qcom stocks is representative of the vast majority of stocks, with the autocorrelation within the noise level already at the minimum delay ( $\tau = 1$ ). This is in good agreement with what we know about financial time series (no dependencies for daily frequency, the typical time scale for autocorrelation being of the order of minutes). We also found that the first point, the autocorrelation between two successive days, is small but negative for most stocks in our data set. In the rightmost panel (MJN stock) we observe a rare dynamic, where the one-lagged autocorrelation is breaching the noise level and then rapidly oscillates to zero.

As clear from the figure, our model reproduces well the observed autocorrelation in all these different cases, and gives a single mathematical explanation for both the exponentially decaying (from positive one-lagged autocorrelation) and the oscillating (from negative one-lagged autocorrelation) behaviour. Moreover, the generic feature of the one-dimensional Ising model, i.e. the absence of a phase transition characterized by a diverging length (here, time) scale [46], explains why in real-world time series the memory is always found to be short-ranged.

#### 4.4. Comparing the three models on empirical financial time series

As we illustrated in section 3.3 in the general case, once we have more than one model for the same data  $\mathbf{X}^*$ , we can use the AIC weights to rank all models in terms of the achieved trade-off between accuracy (good fit to the data) and parsimony (small number of parameters). The AIC

<sup>1</sup> For a set of  $T$  independent and identically distributed pairs of random variables  $\{x_i, y_i\}_{i=1}^T$ , the Pearson correlation coefficient  $\rho_{x,y}$  is distributed around zero, but in a non-Gaussian way. However, the quantity  $\phi_{x,y} \equiv \text{artanh}(\rho_{x,y})$ , known as the *Fisher transformation*, is normally distributed around zero, with standard deviation  $\sigma = (T - 3)^{-1/2}$ . The interval  $-\sigma < \phi_{x,y} < +\sigma$ , representing a 95% confidence interval for  $\phi_{x,y}$ , can then be mapped back to the interval  $-\tanh(\sigma) < \rho_{x,y} < +\tanh(\sigma)$  to obtain a 95% confidence interval for  $\rho_{x,y}$  around zero.



**Figure 6.** Measured AIC weights for the three models (black: uniform random walk; orange: biased random walk; green: one-lagged model) calculated for three different S&P500 stocks, as a function of the time horizon  $T$ . The latter represents the number of months elapsed backwards from October 2011: for all stocks, all time series used to calculate the AIC weights have the same endpoint  $T_0 = 31$  October 2011, and a variable startpoint  $T_0 - T$ .

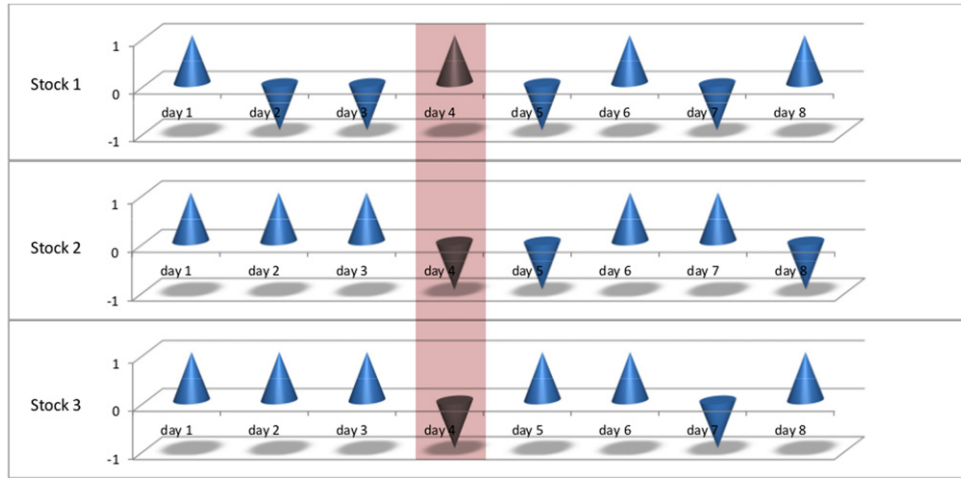
weight  $w_k$  of a specific model  $k$  represents the probability that the model is the ‘best’ one among the candidate models.

We applied this procedure to the three models discussed so far (uniform random walk, biased random walk, one-lagged model). As an example, in figure 6 we show the values of the AIC weights for three different S&P500 stocks. We can see that the performance of the models is wildly fluctuating and different across stocks. This suggests that the informativeness of the measured properties is dependent on different factors, which are not entirely revealed to us. However, it is clear that in all cases the time horizon  $T$  plays a key role in the performance of the models. This means that the outcome depends on how many time steps are included in the analysis. For instance, we see that in some cases (Citigroup Inc. stock) the small  $T$  regime is oscillatory, while the large  $T$  regime appears to set a preference for a definite model. In other cases (United Health Group), the three models alternate over quite long periods of time. Most likely, this very irregular behaviour is due to the strong non-stationarity of financial markets: extending the analysis over longer time horizons does not necessarily improve the statistics, because for large  $T$  the underlying price (and return) distributions change in an uncontrolled way.

We stress again that the AIC weight indicates which property, among the constraints defining all models, can better characterize the stock, given the observed data. In other words, it highlights the *measured property* that is most informative about the original data. Despite the fact that the models considered so far are extremely simplified (and are by no means intended to be accurate models of financial time series), this approach can always identify, in relative terms, the most useful empirical quantity characterizing an observed time series.

## 5. Single cross-sections of multiple time series

In the previous section we considered models for single time series, where  $N = 1$  and  $T$  is large. Here we consider, as a second specification of our general formalism, the somewhat ‘opposite’ case of single cross-sections of  $N$  multiple time series, which represent a daily snapshot of the market dynamics. For clarity, figure 7 portrays a single cross-section of a set of multiple time series. In this case,  $T = 1$  and we assume  $N \gg 1$ . So the matrix  $\mathbf{X}$  has dimensions  $N \times 1$ , i.e. it is



**Figure 7.** An example of a cross-section (highlighted in red) of a set of  $N = 3$  multiple time series. Each cross-section is a  $N \times 1$  matrix (column vector) where each element is the instantaneous binary return of a different stock. For example, the highlighted cross-section is the vector for day  $t = 4$ .

an  $N$ -dimensional column vector. The entries of a cross-section  $\mathbf{X}$  will be denoted by  $x_i$ , where  $1 \leq i \leq N$ , each representing the daily increment of a different asset.

Using again the symbol  $\{ \cdot \}$  to denote an average over stocks (as in section 2.2), we now define the average increment (first moment) of  $\mathbf{X}$  as

$$M_1(\mathbf{X}) \equiv \{x_i\} = \frac{1}{N} \sum_{i=1}^N x_i \quad (51)$$

and the second moment as

$$M_2(\mathbf{X}) \equiv \{x_i^2\} = \frac{1}{N} \sum_{i=1}^N x_i^2 = 1. \quad (52)$$

Therefore the sample variance is

$$M_2(\mathbf{X}) - M_1^2(\mathbf{X}) = 1 - \{x_i\}^2. \quad (53)$$

We also define the total ‘coupling’ between stocks (for a specific cross-section  $\mathbf{X}$ ) as

$$D(\mathbf{X}) \equiv \sum_{i < j} x_i x_j = \{x_i x_j\} \frac{N(N-1)}{2}, \quad (54)$$

where now, as in equation (6),  $\{ \cdot \}$  denotes an average over all pairs of stocks.

In what follows, we will consider various models for single cross-sections. The main difference with respect to the models of single time series considered in section 4 is that the interaction between time steps for a given stock is now replaced by the interaction between different stocks for a given time step. As is well known, in real financial markets the interactions among stocks (as measured, e.g., via cross-correlations) are much stronger than inter-temporal autocorrelations. This makes the cross-sectional properties significantly different from those of the dynamics of single time series, once inter-stock interactions are enforced in the model. Yet, in simple models without interaction, we recover similar expected properties.

### 5.1. Uniform random walk

As in section 4.1, we first consider a trivial model without constraints (see the appendix), defined by the Hamiltonian

$$H(\mathbf{X}) = 0. \quad (55)$$

The probability of the occurrence of a cross-section  $\mathbf{X}$  is completely uniform over the ensemble of all binary cross-sections of  $N$  stocks. Again, this ‘gas of non-interacting spins in vacuum’ model results in a uniform random walk, where all the  $N$  elements of  $\mathbf{X}$  are mutually independent and identically distributed.

In the financial setting, this model assumes that all stocks fluctuate independently of each other (where the ‘fluctuations’ are intended as ensemble ones, since we are now considering a single cross-section), and under the effect of no common factor. Each stock has zero expected value

$$\langle x_i \rangle = 0 \quad (56)$$

and maximum variance

$$\text{Var}[x_i] \equiv \langle x_i^2 \rangle - \langle x_i \rangle^2 = 1. \quad (57)$$

In section 5.4, we will compare the performance of this trivial benchmark to that of the other models we are about to introduce. To this end, the AIC value can be calculated from equation (22) choosing  $n_k = 0$  and using the (constant) likelihood given by equation (B.3) in the appendix.

### 5.2. Biased random walk

In this model, which is analogous to that defined in section 4.2, the constraint is chosen as the total daily increment of the cross-section  $\mathbf{X}$ :

$$C(\mathbf{X}) = N \cdot M_1(\mathbf{X}) = N \cdot \{x_i\}, \quad (58)$$

where  $M_1(\mathbf{X})$  is defined by equation (51). The Hamiltonian is then

$$H(\mathbf{X}, \theta) = \theta \sum_{i=1}^N x_i. \quad (59)$$

Similarly to its counterpart for single time series, this is a model of non-interacting spins under the effect of a common external field, and leads to a biased random walk (see the appendix). The financial interpretation is however different: in this model, all stocks are assumed to fluctuate (again, in an ‘ensemble’ sense) under the effect of a common market-wide factor, but are conditionally independent of each other, given the market-wide factor itself. In the econophysics literature, the overall tendency of all stocks to move together is generally referred to as the ‘market mode’ [2]. When applied to the data, this extremely simple model interprets the observed market mode as the consequence of an external factor (e.g. news), and not of direct interactions among stocks.

The probability  $P_i(x | \theta)$  of a given increment  $x = \pm 1$  for stock  $i$  is

$$P_i(x | \theta) = \frac{e^{-\theta x}}{e^{-\theta} + e^{+\theta}}, \quad (60)$$

the expected value of the  $i$ th increment  $x_i$  is

$$\langle x_i \rangle_\theta = -\tanh \theta, \quad (61)$$

and the variance is

$$\text{Var}[x_i] = 1 - \tanh^2 \theta. \quad (62)$$

The maximum likelihood condition (16), fixing the value  $\theta^*$  of the parameter  $\theta$  given a real cross-section  $\mathbf{X}^*$ , leads to

$$\theta^* = -\frac{1}{2} \ln \left[ \frac{1 + \{x_i^*\}}{1 - \{x_i^*\}} \right], \quad (63)$$

where  $\{x_i^*\}$  is the measured average increment of the observed cross-section  $\mathbf{X}^*$ . We will apply this model to real financial data in sections 5.4 and 6. The AIC of the model is given by equation (22) where  $n_k = 1$  and where the maximized likelihood is given by  $P(\mathbf{X}^* | \theta^*)$ , with  $P(\mathbf{X} | \theta)$  given by equation (B.10) (see the appendix).

### 5.3. Mean field model

We now consider a more complex model, with interactions among *all* stocks, which is suitable for financial cross-sections. Besides the constraint on the total increment, we enforce an additional constraint on the average coupling between stocks. The resulting two-dimensional constraint can be written as

$$\vec{C}(\mathbf{X}) = \begin{pmatrix} C_1(\mathbf{X}) \\ C_2(\mathbf{X}) \end{pmatrix} = \begin{pmatrix} N \cdot M_1(\mathbf{X}) \\ D(\mathbf{X}) \end{pmatrix}, \quad (64)$$

where  $M_1(\mathbf{X})$  is given by equation (51) and  $D(\mathbf{X})$  by equation (54). If we write the corresponding Lagrange multiplier as

$$\vec{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = - \begin{pmatrix} h \\ J \end{pmatrix} \quad (65)$$

then the Hamiltonian reads

$$H(\mathbf{X}, h, J) = -h \sum_{i=1}^N x_i - J \sum_{i < j} x_i x_j. \quad (66)$$

Like the one-lagged model for single time series (see section 4.3), this model is formally analogous to an Ising model of interacting spins under the influence of an external ‘magnetic’ field (here denoted by  $h$ ). However, the big difference is that, whereas in the one-lagged model each increment  $x(t)$  interacts *only with the next temporal increment*  $x(t + 1)$  of the same stock, here each increment  $x_i$  interacts *with all the other increments*  $x_j$  of the same cross-section  $\mathbf{X}$ , i.e. with all other stocks in the market. As a model of spin systems, the above model is generally known as the mean-field Ising model [46]. In the appendix we provide the analytical solution of the model, adapted to our setting.

In the financial setting, this model allows us to separately consider the effects of the external field, i.e. a common factor affecting all stocks in the market, from those of the average



interaction among all stocks. This market-wide interaction can also cause all stocks to correlate, but has the different interpretation of a collective effect, i.e. the tendency of stock increments to ‘align’ with each other as a result of direct interactions, rather than of a common influence. This is a sort of ‘herd effect’ at the coarse-grained level of attractive ( $J > 0$ ) inter-stock interactions. So, the model can generate the ‘market mode’ either as the result of a common external influence such as news (in which case all stocks are still conditionally independent given the common factor), or as a collective effect due to mutual interactions (in which case all stocks are conditionally dependent given the common factor).

While the model can in principle simulate synthetic time series under a combination of the above two effects by varying the two parameters  $h$  and  $J$  independently, a problem arises when it is fitted to the data. The mathematical root of the problem is the well known fact that  $H(\mathbf{X}, h, J)$  can be rewritten as a linear combination of  $M_1(\mathbf{X})$  and  $M_1^2(\mathbf{X})$ . As we show in the appendix, this implies that, when the maximum likelihood principle is used to fit the model to the data  $\mathbf{X}^*$ , the variance of  $M_1(\mathbf{X})$  becomes zero. In other words, the model degenerates to one where  $M_1(\mathbf{X})$  is no longer a random variable. This also implies that the two equations fixing the values of the parameters  $J^*$  and  $h^*$  become identical (see the appendix). Therefore it is no longer possible to uniquely fix the values of both parameters, and the problem is over-constrained. For this reason, we need to eliminate one parameter and consider the model only in the two extreme cases  $h = 0$  and  $J = 0$ . These two cases can be treated separately.

The case  $J = 0$  coincides with the biased random walk model already considered in section 5.2, where  $\theta = -h$ . Using equation (63), we therefore specify this model using the two parameter values

$$h^* = \frac{1}{2} \ln \left[ \frac{1 + \{x_i^*\}}{1 - \{x_i^*\}} \right], \quad J^* = 0, \quad (67)$$

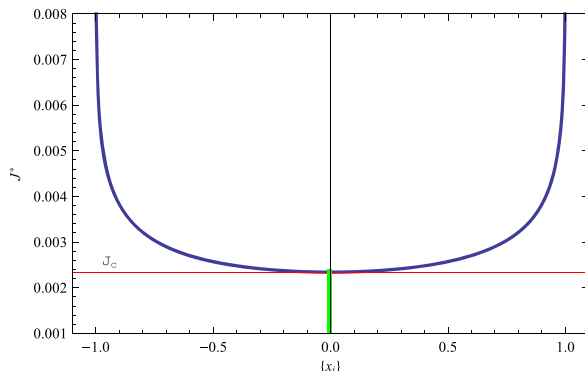
where  $\{x_i^*\}$  is the observed average increment of the empirical cross-section  $\mathbf{X}^*$ . This model interprets the market mode as arising *only* from a common external factor.

The case  $h = 0$  leads us instead to a novel model where the market mode is interpreted *only* as a collective effect arising from inter-stock interactions. Using the analytical results reported in the appendix, and in particular equation (B.35), we find that the parameter values are in this case

$$h^* = 0, \quad J^* = \frac{1}{2\{x_i^*\}(N-1)} \ln \left[ \frac{1 + \{x_i^*\}}{1 - \{x_i^*\}} \right]. \quad (68)$$

In what follows, when using the ‘mean-field’ model, we will always refer to the parameter specification defined by (68). The other specification, equation (67), will instead still be denoted as the ‘biased random walk’ model.

In figure 8 we plot the value of  $J^*$  as a function of  $\{x_i^*\}$ , as defined by equation (68). We note however that equation (68) is undefined for  $\{x_i^*\} = \pm 1$  and  $\{x_i^*\} = 0$ . The breakdown for  $\{x_i^*\} = \pm 1$  simply means that, in order to align *all* returns (in either direction),  $J^*$  should diverge to  $+\infty$ . The breakdown for  $\{x_i^*\} = 0$  is instead more profound. For infinitesimal (both positive and negative) values of  $\{x_i^*\}$ ,  $J^*$  admits the finite limit



**Figure 8.** The value of the fitted parameter  $J^*$  as a function of the measured average binary return  $\{x_i^*\}$  (blue curve) for a group of  $N = 428$  stocks (as in our S&P sample). The curve shows a one-to-one relationship for  $\{x_i^*\} \neq 0$ . While  $\lim_{\{x_i^*\} \rightarrow 0} J^* = J_c \equiv (N - 1)^{-1}$ , for  $\{x_i^*\} = 0$  the value of  $J^*$  is actually indeterminate, as there is an infinity of values of  $J^*$  (namely all values  $-\infty < J^* \leq J_c$ , see vertical green line) that are possible solutions of the model. The value of  $J_c$  is indicated by the horizontal red line.

$$\lim_{\{x_i^*\} \rightarrow 0^+} J^* = \lim_{\{x_i^*\} \rightarrow 0^-} J^* = \frac{1}{N - 1}. \quad (69)$$

However, at the very point  $\{x_i^*\} = 0$ ,  $J^*$  is actually indeterminate.

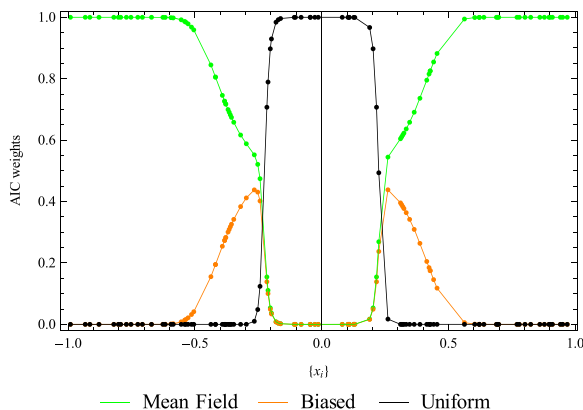
The above effect is due to the well-known phase transition of the mean-field Ising model. In the traditional physical setting, the phase transition occurs at a critical temperature (here reabsorbed in the value of the parameters  $h$  and  $J$ ). When  $h = 0$ , the critical value is obtained by setting  $(N - 1)J = 1$ , because for  $(N - 1)J < 1$  equation (B.33) (see the appendix) has the single solution  $\langle M_1 \rangle = 0$ , corresponding to a phase with no macroscopic magnetization, while for  $(N - 1)J > 1$  there are three solutions, one of which is still  $\langle M_1 \rangle = 0$  (which is now unstable) and the other two being the stable solutions  $\langle M_1 \rangle = \pm m$  (corresponding to the onset of a macroscopic magnetization  $|m| > 0$  where most spins point in the same direction). In our financial setting, since the magnetization is fixed by the data through the relation  $\langle M_1 \rangle = \{x_i^*\}$ , the condition  $(N - 1)J^* = 1$  implies that the phase transition occurs at the critical value

$$J_c = \frac{1}{N - 1} \quad (70)$$

of the control parameter  $J^*$ . We can therefore rewrite equation (69) as

$$\lim_{\{x_i^*\} \rightarrow 0} J^* = J_c. \quad (71)$$

For  $J^* > J_c$  we get a ‘magnetized’ phase where most stock prices move in the same direction (aligned returns), while for  $J^* < J_c$  we get a non-magnetized phase where there is no collective alignment of stock increments, and  $\{x_i^*\} = 0$ . We therefore conclude that the reason why the value of  $J^*$  is indeterminate for  $\{x_i^*\} = 0$  is because there is an infinity of values of  $J^*$  (namely all values  $-\infty < J^* \leq J_c$ ) that are possible solutions of the model.



**Figure 9.** The calculated AIC weights of the three cross-sectional models (uniform random walk, biased random walk, mean field model) as a function of the measured average daily binary return  $\{x_i^*\}$ , for  $N = 428$  S&P500 stocks, each studied for 100 days of trade.

It should be noted that the case  $\{x_i^*\} = 0$  is never practically encountered in reality, since the empirical  $\{x_i^*\}$  can be arbitrarily small, but is generally not really zero. While this ‘protects’ the model from the indeterminacy discussed above, it raises another problem of arbitrariness, which can however be solved very effectively using the information-theoretic criteria that we have introduced in section 3.3. The problem is that the mean-field model will always interpret even the tiniest empirical deviations from  $\{x_i^*\} = 0$  as the result of direct interactions among stocks, and attach a value  $J^* > 0$  to this interpretation. This will also apply to, e.g., most realizations of a purely uniform random walk: even if for such a model one knows that the theoretical expected return is zero, most realizations will be such that  $\{x_i^*\}$  is small but non-zero. So the only phase of the mean-field model that can be explored is the ‘magnetized’ phase dominated by collective effects. This implies that even a pure effect of noise will be interpreted as the presence of interactions. However, this problem will be solved in the next section, where we show that an information-theoretic comparison between the mean-field model, the uniform random walk, and the biased random walk is able to discriminate the most parsimonious model, thus allowing us to trust the mean-field model only when  $\{x_i^*\}$  is distant enough from zero.

#### 5.4. Comparing the three models on empirical financial cross-sections

We can now combine the three models together and use the AIC weights (see section 3.3) to determine which model achieves the optimal trade-off between accuracy and parsimony. This will immediately provide us with an indication of whether the observed market mode, as reflected in the empirical aggregate increment  $\{x_i^*\}$ , should be interpreted, e.g., as a common exogenous factor, as a collective endogeneous effect, or even only as the sheer outcome of chance.

The fact that the likelihoods of the biased random walk and the mean-field model depend only on  $\{x_i^*\}$  and  $N$ , plus the fact that the likelihood of the uniform random walk is constant, allows us to obtain the AIC values for the three models as functions of  $\{x_i^*\}$  and  $N$  only. In figure 9 we show the calculated AIC weights of the three models as a function of the observed value  $\{x_i^*\}$ , for  $N = 428$  S&P500 stocks. Each point represents a different cross-section, i.e. a

different day of trade, for a total of 100 randomly sampled days. It is important to note that the empirical value of the average increment only determines which point(s) of the curves are actually visited, but the curves themselves are universal.

The figure reveals a remarkable fact, namely the presence of three distinct regimes in the behavior of the group of stocks. For  $0 \leq |\{x_i^*\}| \lesssim 0.2$ , we find that the best performing model is the uniform random walk, which displays an AIC weight practically equal to one (indicating that the model is almost surely the best one among the three models considered, see section 3.3). This means that, in this ‘noisy’ regime, the most parsimonious explanation of the market mode, as reflected in the measured value of  $\{x_i^*\}$ , is that of a pure outcome of chance.

For  $0.2 \lesssim |\{x_i^*\}| \lesssim 0.5$ , we find that the uniform random walk is almost surely *not* the best model, while the biased random walk and mean field models are competing. We observe an almost equal performance of the two models for  $|\{x_i^*\}| \approx 0.2$ , and an increasing preference for the mean field model as  $|\{x_i^*\}|$  increases towards 0.5. Despite this preference, we cannot reject the mean field model, meaning that in this ‘mixed’ regime the most likely explanation for the market mode is a combination of exogenous and endogenous effects.

Finally, for  $0.5 \lesssim |\{x_i^*\}| \lesssim 1$ , the mean field model achieves practically unit probability to be the best model. In this ‘endogenous’ regime, the most likely explanation for the market model is uniquely in terms of a collective effect of direct influence among stocks.

We can summarize the above findings as follows:

$$\left\{ \begin{array}{ll} \text{Uncoordinated (noisy) regime:} & 0 \leq |\{x_i^*\}| \lesssim 0.2 \\ \text{Mixed (endogenous + exogenous) regime:} & 0.2 \lesssim |\{x_i^*\}| \lesssim 0.5 \\ \text{Coordinated (endogenous) regime:} & 0.5 \lesssim |\{x_i^*\}| \leq 1, \end{array} \right.$$

where we recall that the values of  $|\{x_i^*\}|$  delimiting the various regimes have been calculated for  $N = 428$ .

While the qualitative finding that larger values of  $|\{x_i^*\}|$  are better explained in terms of collective effects might appear intuitive, the possibility to quantitatively identify the value  $|\{x_i^*\}| \approx 0.5$  above which this intuition is fully supported by statistical evidence is a non-obvious output of the above approach. The same consideration applies to the identification of the other two regimes, and of a mixed phase where there is not enough statistical evidence in favour of a single interpretation of the market mode. Moreover, the fact that the mean field model starts being statistically significant only for  $|\{x_i^*\}| \gtrsim 0.2$  solves the aforementioned problem of an otherwise problematic interpretation of even tiny values of  $|\{x_i^*\}|$  as the result of inter-stock interactions. The AIC analysis shows that, for values below 0.2, one should not trust the mean field model, and consequently the value  $J^* > 0$  that the model itself indicates. When  $|\{x_i^*\}| \lesssim 0.2$ , the best model is actually the uniform random walk, which effectively corresponds to  $J^* = 0$ . This is a highly non-trivial result.

## 6. Ensembles of matrices of multiple time series

In this section, as our third and final specification of the abstract formalism introduced in section 3, we extend the previous results to the general case where the observed data is a full  $N \times T$  matrix  $\mathbf{X}^*$  representing a set of multiple binary time series for  $N$  stocks, each extending

over  $T$  timesteps. We recall that the entries of a generic such matrix  $\mathbf{X}$  are denoted by  $x_i(t)$ , where  $i$  labels the stock and  $t$  labels the time step. We assume that  $N$  and  $T$  are both large, i.e.  $N \gg 1$  and  $T \gg 1$ . Before introducing an explicit model, we need to make some important considerations.

We had already anticipated that the purpose of the models introduced in the previous sections was not that of introducing realistic models of financial time series. For instance, it is well known that the simple stochastic processes considered in section 4 are far too simple to reproduce some key stylized facts observed in real financial time series, such as volatility clustering [47, 48] or bursty behavior [49]. Moreover, being entirely binary, the above examples cannot address other well established properties characterizing the amplitude of fluctuations, e.g. the ‘fat’ (power-law) tails of the empirical distributions of price returns.

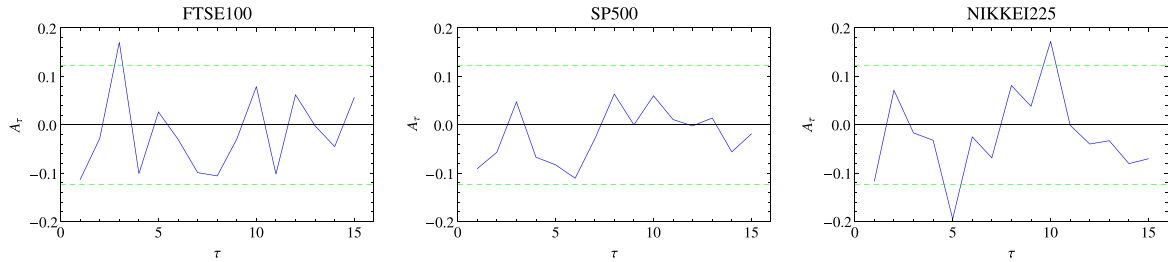
Nonetheless, there is a simple argument that legitimates us to use a proper extension of the above modelling approach, especially that introduced in section 5, provided that we adequately calibrate such extension on the observed set of multiple time series. The argument is basically the realization that we can properly model the binary signature of a time series, using temporal iterations of even the simplistic models we have introduced in section 5, if we assume that some aggregated information measured on the original ‘weighted’ time series  $r_i(t)$  ( $1 \leq i \leq N$ ) can be used as a proxy of the driving factor defining the model itself. We will show that this simple assumption is actually verified in the data. In particular, we will show that a sequence of temporal iterations of the biased random walk model, which assumes that the binary time series is driven by an ‘external’ field, can be ‘bootstrapped’ on the real data by assuming that the field can be replaced by a function of the (endogenous) observed aggregate increment of the original weighted time series, i.e. the empirical value  $\{r_i^*\}$  of the quantity  $\{r_i\}$  defined in equation (3). In such a way, we do not need a model generating a realistic dynamics of  $\{r_i\}$  (or of the individual stock-specific increments) in order to model the behaviour of  $\{x_i\}$ , because the time series of  $\{r_i\}$  is taken from the data.

As a result, we will obtain an accurate model for the dynamics of the aggregate binary increment  $\{x_i(t)\}$ , given the observed dynamics of  $\{r_i(t)\}$ . This model will reproduce with great accuracy, and mathematically characterize, the empirical nonlinear relation between these two quantities that we have illustrated in section 2.2. We will finally test the temporal robustness and predictive power of the model, and conclude with a discussion of the relatedness of our approach and more traditional ‘factor models’ in finance.

### 6.1. Temporal dependencies among cross-sections

In order to execute the above plan, we first analyze the correlations between single cross-sections of the market. We need this preliminary analysis in order to determine whether the temporal extension of the models defined in section 5 should incorporate dependencies among different snapshots.

Based on extensive financial literature, we expect no correlation (on a daily frequency) among the returns of different cross-sections. However, most analyses focus on the auto-correlation of *individual* stocks, based on their *weighted* returns. So, to check our hypothesis we perform an explicit analysis of the temporal auto-correlation of the observed time series of the *aggregate, binary* return  $\{x_i^*(t)\}$ . This analysis is shown in figure 10 for the three indices, using daily data for the year 2006. We confirm that the observed autocorrelation is not statistically significant, since (apart for a few points) it lies within the range of random noise (calculated by



**Figure 10.** The measured autocorrelation of the average binary daily return  $\{x_i^*(t)\}$  for the three indices in year 2006. The green lines represent the noise level, calculated as  $\pm 2$  standard deviations of the Fisher-transformed autocorrelation.

imposing a threshold of two standard deviations on the Fisher-transformed autocorrelation). This type of uncorrelated dynamics is observed throughout our dataset. This means that, in line with other analyses of autocorrelation, the memory of the aggregate binary return of real markets, if any, is much shorter than a day.

Going back to the result illustrated in figure 9, we can then conclude that there is no significant correlation in the trajectories of the daily points populating the curves. In other words, given the knowledge of the position of the market in the AIC curves in a given day, we cannot predict where the market will move the next day, even if of course we know that it will move to another point in the curves themselves.

## 6.2. Reproducing the observed binary/non-binary relationships

The previous result sets the stage for our next step, where we consider an explicit extension of the models considered in section 5 to an ensemble of multiple time series, as introduced in section 3 in the general case. The absence of autocorrelation implies that we can define the Hamiltonian of the full  $N \times T$  matrix  $\mathbf{X}$  as a sum of  $T$  non-interacting Hamiltonians, each describing a single cross-section of  $N$  stocks.

Next, we need to choose the model to extend. We want the final model to establish (among other things) an expected relationship between the binary and the weighted aggregate returns, so that we can test this prediction against the empirical relationships illustrated in section 2.2. This implies that we need to input the measured weighted return  $\{r_i^*\}$  as a driving parameter of the binary model. Among the three models, only the biased random walk and the mean field model have parameters that can be related to  $\{r_i^*\}$ . In section 5 we treated those models as giving competing interpretations of the market model in terms of exogenous and endogenous effects, respectively. However, it should be noted that this is no longer possible as soon as the parameters of these models are made dependent on the observed return. For instance, if we assume that the parameter  $\theta$  of the biased random walk depends on  $\{r_i^*\}$  (which is a property of the data), we can no longer interpret  $\theta$  as an external field, since it has been somehow ‘endogenized’. Determining whether  $\theta$  can be interpreted as endogenous or exogenous is now entirely dependent on whether  $\{r_i^*\}$  itself can be interpreted as endogenous or exogenous. This tautology does not prevent us from determining a relationship between  $\{r_i^*\}$  and  $\{x_i^*\}$  in their full range of variation, because such a relationship is independent on the optimal (endogenous or exogenous) interpretation of both quantities.



We also note that the choice of the model to calibrate on  $\{r_i\}$  is now completely independent of the relative performance of the various models that we have determined in the case of free parameters, including their AIC weights shown in figure 9. Indeed, apart from an initial calibration, the parameters will no longer be fitted using the ML principle, making the AIC analysis no longer appropriate. In other words, ranking the ‘free’ models and endogenizing their parameters are two completely different problems. In particular, the low AIC weight of the biased random walk throughout most of figure 9 does not impede us from using this model in our next analysis. We will indeed ‘bootstrap’ the biased random walk on the real data, by looking for a relationship between  $\{r_i\}$  and the parameter  $\theta$ . We prefer this model over the mean field one because, while it is natural to think of (a function of)  $\{r_i\}$  as a proxy of the ‘field’  $\theta$  affecting the market in the biased random walk model (notably,  $\{r_i\}$  has a definition similar to that of a market index), it is less natural to think of the same quantity as a proxy of the inter-stock interaction  $J$  in the mean field model (although, as we said before, this would be technically possible).

Combining all the above considerations, we finally generalize the biased random walk model defined by equation (59) to the matrix case as follows:

$$H(\mathbf{X}, \vec{\theta}) = \sum_{t=1}^T \theta(t) \sum_{i=1}^N x_i(t), \quad (72)$$

where  $\vec{\theta}$  is a  $T$ -dimensional vector with entries  $\theta(t)$ . Note that, while the models we introduced in section 4 have time-independent parameters and therefore correspond to time series at statistical equilibrium (for example a model with constant volatility), we are now considering more general models with time-dependent parameters. Relating  $\theta(t)$  to  $\{r_i(t)\}$  will allow us to incorporate any observed degree of non-stationarity of the data into the model itself.

As a preliminary calibration, we now look for an empirical relation between  $\{r_i(t)\}$  and  $\theta(t)$ . To this end, we first treat the latter as a free parameter and look for the optimal value  $\theta^*(t)$  maximizing the likelihood of the observed binary time series  $\mathbf{X}^*$ . Since the Hamiltonians for different timesteps are non-interacting, it is easy to show that  $\theta^*(t)$  is given again by equation (63) where  $\{x_i^*\}$  is replaced by  $\{x_i^*(t)\}$ :

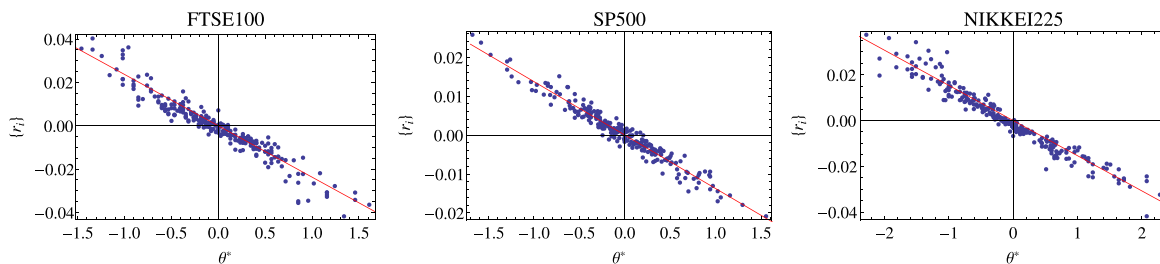
$$\theta^*(t) = -\frac{1}{2} \ln \left[ \frac{1 + \{x_i^*(t)\}}{1 - \{x_i^*(t)\}} \right]. \quad (73)$$

In figure 11 we compare the resulting value of  $\theta^*(t)$  with the corresponding observed weighted return  $\{r_i^*(t)\}$ , for the three indices separately. Each point in the plot corresponds to a different day, and we considered 250 days (approximately one year) for each index. We find a strong linear relation between the two quantities. This relation can be fitted by the one-parameter curve

$$\{r_i^*(t)\} = -c \cdot \theta^*(t), \quad (74)$$

where  $c > 0$ . This finding is very important. It confirms that the parameter  $\theta^*(t)$ , defined through equation (72) as a time-varying ‘field’ driving the observed binary increment  $\{x_i^*(t)\}$  with maximum likelihood, is an excellent proxy for the observed non-binary ‘market index’  $\{r_i^*(t)\}$ . This result holds up to a negative factor  $c$  which, on the time scale considered, is constant for each market (in section 6.3 we will provide a more detailed analysis of the stability





**Figure 11.** The most likely value of the driving field  $\theta^*(t)$  calculated by applying the biased random walk model to the projected binary signature of day  $t$ , compared with the measured average weighted return  $\{r_i^*(t)\}$  of the same day, for 250 trading days (approximately one year) in the FTSE100 (left), S&P500 (center) and NIKKEI225 (right) in various years (2003, 2007, and 2004 respectively). We also show the linear fit  $\{r_i(t)\} = -c\theta^*(t)$  with  $c > 0$ .

of  $c$  over different time scales). Since  $\{r_i^*(t)\}$  is a property measured on the stock increments themselves, it reflects both external influences and internal dependencies. Therefore  $\theta^*(t)$  cannot be (entirely) interpreted as an external field. This confirms our interpretation of the biased random walk as a model agnostic to the (endogenous or exogenous) nature of the driving field in the present setting.

Combining equations (73) and (74), we finally obtain a mathematical expression for the expected relationship between  $\{r_i^*\}$  and  $\{x_i^*\}$  in our model:

$$\{r_i^*(t)\} = \frac{c}{2} \ln \left[ \frac{1 + \{x_i^*(t)\}}{1 - \{x_i^*(t)\}} \right] = c \cdot \operatorname{artanh} \{x_i^*(t)\}. \quad (75)$$

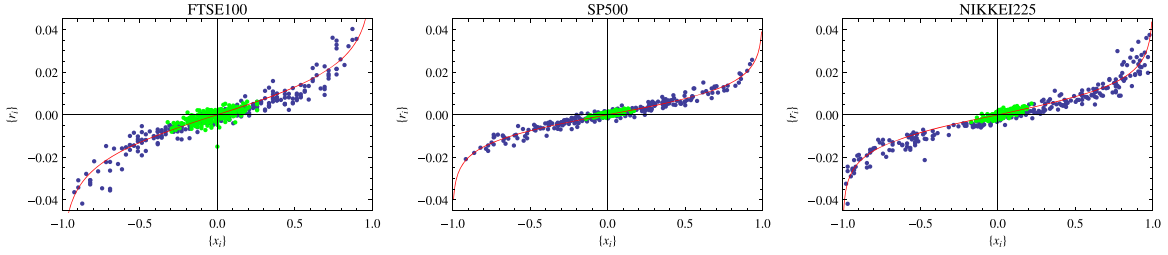
Inverting, we have

$$\{x_i^*(t)\} = \tanh \frac{\{r_i^*(t)\}}{c}. \quad (76)$$

We can now test the above expressions against the data shown previously in figure 2. In that figure, we already showed that the observed relationship between  $\{r_i^*\}$  and  $\{x_i^*\}$  can be fitted very well by a curve of the form given by equation (75). We have just provided a theoretical justification for the otherwise arbitrary use of such expression. Moreover, now we can fit the value of  $c$  using equation (74), which is independent of equation (75). Once we obtain  $c$  in this way, we can use equation (75) to predict  $\{r_i^*(t)\}$  given  $\{x_i^*(t)\}$ , or *vice versa*, without fitting any parameter. In figure 12 we show the result of this operation. We confirm that the prediction of our model matches the empirical relationship very well.

We also consider a null model where we randomly shuffle the increments of each of the  $N$  time series independently. This results in a set of randomized time series, with elements  $r_i'(t)$ , where the total increment  $\sum_{t=1}^T r_i'(t)$  for each stock is preserved, but the returns of all stocks in a given day are uncorrelated. From  $r_i'(t)$ , we obtain the binary signature  $x_i'(t)$  as for the real data. As shown in figure 12, this randomized benchmark overlaps with the empirical trend only in a very narrow, linear regime. We will now try to understand this result.

The reason why the shuffled data result in a linear trend is the following. For each value of  $\{x_i'\}$ , there is a definite number  $N_{\text{up}}$  of ‘up’ stocks and a definite number  $N_{\text{down}} = N - N_{\text{up}}$  of



**Figure 12.** Nonlinear relationship between the average daily increment (weighted return) and the average daily sign (binary return) over all stocks in the FTSE100 (left), S&P500 (center) and NIKKEI225 (right) in various years (2003, 2007, and 2004, respectively). Here each point corresponds to one day in the time interval of 250 trading days (approximately one year). The red curve is our non-parametric prediction based on the fit shown in figure 11, and the green points are the same properties measured on the shuffled data.

‘down’ stocks, according to the relation

$$\{x'_i\} = \frac{N_{\text{up}} - N_{\text{down}}}{N} = \frac{2N_{\text{up}} - N}{N} = 2\frac{N_{\text{up}}}{N} - 1. \quad (77)$$

Conditional on the above value of  $\{x'_i\}$ , the expected value of  $\{r'_i\}$  (over multiple shufflings) is

$$\langle \{r'_i\} \rangle = \frac{r_+^* N_{\text{up}} + r_-^* N_{\text{down}}}{N} \approx \frac{r_+^* N_{\text{up}} - r_+^* N_{\text{down}}}{N} = r_+^* \left[ 2\frac{N_{\text{up}}}{N} - 1 \right] = r_+^* \{x'_i\}, \quad (78)$$

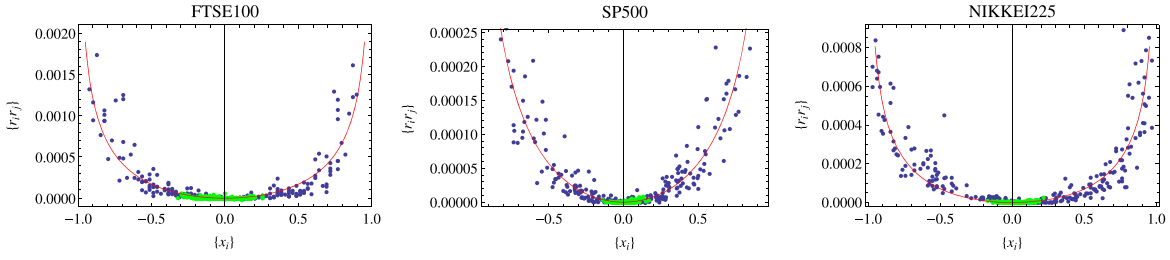
where  $r_+^* > 0$  is the average positive increment (over all  $T$  time steps and all  $N$  time series) and  $r_-^* < 0$  is the average negative increment. Note that both values coincide with the corresponding quantities in the original data, and have been denoted by a star accordingly. Assuming approximately symmetric log-return distributions for each of the  $N$  time series as typically observed, we have set  $r_-^* \approx -r_+^*$ . Given the overlap between real and shuffled data around zero returns in figure 12, we can linearize equation (76) around zero and compare it with equation (78) to get

$$c \approx r_+^*. \quad (79)$$

The above expression suggests that the value of  $c$  strongly depends on the original log-return distribution. Therefore, we expect that the stability of  $c$  is determined by that of  $r_+^*$ . In section 6.3 we will study the stability of  $c$  in more detail.

The above simple argument shows that, for shuffled data, we indeed expect a linear relationship between  $\{r'_i(t)\}$  and  $\{x'_i(t)\}$ . This is a striking difference with respect to real data, where  $\{r_i^*(t)\}$  virtually diverges as  $|\{x_i^*(t)\}|$  approaches one. This ‘divergence’ indicates that, when most stocks are aligned in real markets ( $|\{x_i^*(t)\}| \approx 1$ ), the observed log-returns are much larger than the typical positive increment ( $|\{r_i^*(t)\}| \gg r_+^*$ ). In other words, extreme log-returns are more often observed when stocks are synchronized. This means that there is a strong correlation between the magnitude of log-returns of individual time series and the degree of coordination of all stocks in the market.

While for infinite realizations of the shuffling procedure we would observe equation (78) extending to the full range  $-1 \leq \{x'_i\} \leq +1$ , for finite realizations we observe a much narrower span of values (see figure 12). This is due to the absence of correlations among stocks, resulting



**Figure 13.** Nonlinear relationship between the average daily coupling (weighted coupling) and the average daily sign (binary return) over all stocks in the FTSE100 (left), S&P500 (center) and NIKKEI225 (right) in various years (2003, 2007, and 2004, respectively). Here each point corresponds to one day in the time interval of 250 trading days (approximately one year). The red curve is our non-parametric prediction based on the fit shown in figure 11, and the green points are the same properties measured on the shuffled data.

in significantly lower values of both  $\{r'_i\}$  and  $\{x'_i\}$  with respect to the observed quantities  $\{r_i^*\}$  and  $\{x_i^*\}$ . Interestingly enough, for the S&P500 index the randomized data span the range  $|\{x'_i\}| \lesssim 0.2$ , which coincides precisely with the regime we identified in figure 9 for a completely noisy-driven system with the same number of stocks. This confirms that the AIC analysis correctly pinpoints the boundaries outside which one should expect the observed value  $\{x_i^*\}$  to be inconsistent with a typical realization of  $N$  purely random variables.

The above results also provide an explanation for the second empirical nonlinear relation that we had documented in section 2.2, i.e. the one between  $\{r_i^*(t)r_j^*(t)\}$  and  $\{x_i^*(t)\}$  (see figure 3). In general, we can write  $\{r_i r_j\}$  as

$$\{r_i r_j\} = \frac{1}{N(N-1)} \sum_{i \neq j} r_i r_j = \frac{1}{N(N-1)} \left[ \left( \sum_{i=1}^N r_i \right)^2 - \sum_{i=1}^N r_i^2 \right]. \quad (80)$$

The term  $\sum_{i=1}^N r_i^2$  is of order  $N$ , and vanishes for large markets when divided by  $N(N-1)$ . We are therefore left with

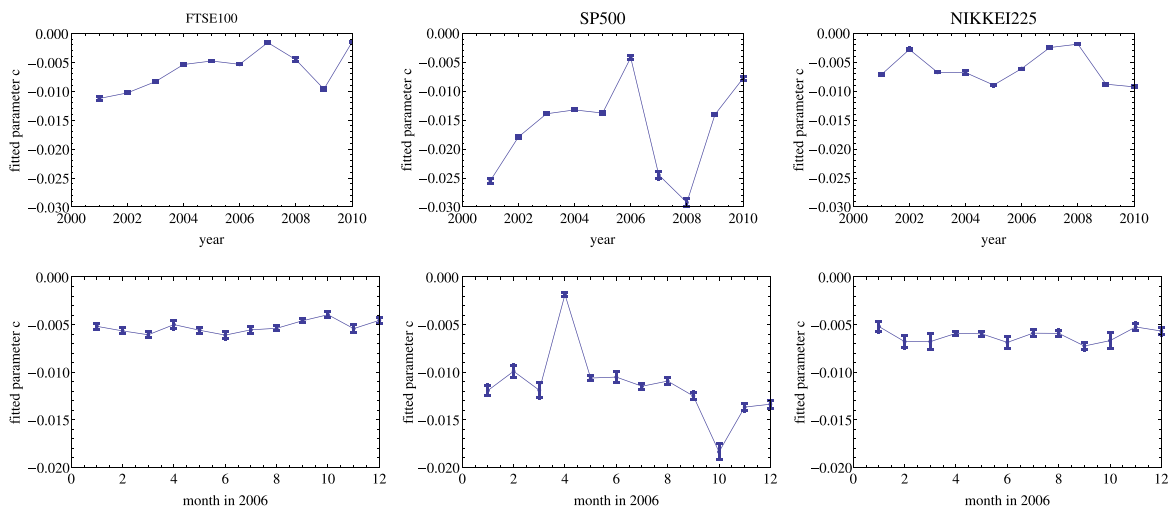
$$\{r_i r_j\} \approx \frac{1}{N(N-1)} \left( \sum_{i=1}^N r_i \right)^2 \approx \{r_i\}^2. \quad (81)$$

Using equation (75), we get

$$\{r_i^*(t)r_j^*(t)\} \approx \{r_i^*(t)\}^2 = c^2 \text{artanh}^2 \{x_i^*(t)\} \quad (82)$$

which theoretically justifies the fitting function we had used in figure 3. Again, rather than fitting that curve on the data, we can use the value of  $c$  determined from the (independent) fit shown in figure 11. This results in the non-parametric plot shown in figure 13. We confirm that, for each of the three indices, we can reproduce the observed relationship very well.

As before, we also show the relationship between  $\{r'_i(t)r'_j(t)\}$  and  $\{x'_i(t)\}$  for randomly shuffled data. The linearity of equation (78) now translates into an expected parabolic relationship:



**Figure 14.** Stability of the parameter  $c$ , fitted as in figure 11 on various yearly (top panels) and monthly (bottom panels) snapshots of the market, for the FTSE100 (left), S&P500 (center) and NIKKEI225 (right).

$$\langle \{r'_i r'_j\} \rangle \approx \{r'_i\}^2 = (r_+^*)^2 \{x'_i\}^2. \quad (83)$$

Again, real data strongly deviate from the above ‘uncorrelated’ parabolic expectation, because extreme events make the empirical coupling  $\{r_i^* r_j^*\}$  virtually ‘diverge’ when stocks are highly synchronized ( $|\{x_i^*\}| \approx 1$ ).

### 6.3. Stability of the parameter $c$

Once we have mathematically characterized the observed nonlinear relations, an unavoidable question arises: in a given market, how stable are those relations? Since  $c$  is the only parameter in the above analysis, the question simply translates into the stability of  $c$ . We have already noted that  $c$  is related to the average positive return  $r_+^*$ , which we expect to be relatively stable. In order to study the stability of  $c$  in more detail, we now consider several yearly and monthly time windows, and explore the time evolution of the fitted parameter for the three indices.

In figure 14 (upper panels) we plot the values of the parameter  $c$  (with error bars) for 11 yearly snapshots (2001–2010). It is clear that there are periods during which the yearly values are relatively stable, and periods when they fluctuate wildly. Thus, in most cases the fitted value of  $c$  in a given year does not allow one to make predictions about the value of  $c$  in the next year.

However, we can also consider a monthly frequency. In the bottom panels of figure 14 we show the result of our analysis, when carried out on the 12 monthly snapshots of year 2006. We choose this particular year because, in the yearly trends shown above, it represents very different points for different markets: the end of a stable period for the FTSE100, an exceptional jump for the S&P500, and the middle of an increasing trend for the NIKKEI225. Despite these differences, we find that in all three markets the monthly dynamics is much more stable than the yearly one. In particular, the trends for FTSE100 and NIKKEI225 are almost constant, and for the S&P500 there are only two deviating points from an otherwise stable trend (despite the large fluctuation that 2006 represents in the yearly trend for this index). This implies that, in most cases, one might even use the monthly value of  $c$  out of sample, in order to predict the future

relationship between  $\{x_i\}$  and  $\{r_i\}$  based on a past observation. We should however stress that the aim of our method is to characterize such relationship, and not to predict it. Indeed, we cannot imagine any situation in which only the binary (or only the non-binary) information is available.

The above results show that there is a trade-off between short and long periods of time. For short (e.g. monthly) periods there are fewer points to calculate  $c$  through a fit of the type shown in figure 11. This explains why the monthly trends in figure 14 have larger error bars than the yearly trends in the same figure. By contrast, for longer (e.g. yearly) periods each individual fit is better, but there are more fluctuations in the temporal evolution of the parameter  $c$ , because the data are less stationary. In general, we expect that in each market, and for a specific period of time, there is a different ‘optimal’ frequency to consider.

#### 6.4. Relation to factor models

We would like to conclude this paper with a discussion of the relationship between some of our findings and the popular *factor models* in the financial literature [3]. As a basic consideration, we stress that factor models can only be applied to the original (non-binary) increments (it is impossible to decompose a binary signal into a non-trivial combination of binary signals), while our models only apply to the binary projections. We should bear this irreducible difference in mind in what follows. However, due to the mapping between binary and non-binary increments that we have documented, we can indeed try to relate the two approaches.

First, let us consider the shuffled (uncorrelated) data, where the original log-returns are randomly permuted within each of the  $N$  time series. It is well known that the total temporal increment (over  $T$  time steps) of any empirical time series of price increments is generally close to zero (due to market efficiency), and that the distribution of log-returns is mostly symmetric around this value. This is especially true if each of the  $N$  original time series has been separately standardized, i.e. the  $i$ th temporal average has been subtracted from each increment of the  $i$ th time series, and the result has been divided by the  $i$ th standard deviation. In such a case, the  $N$  log-return distributions become also very similar to each other, because their support is the same and their values are comparable. This means that, after the shuffling, the time series are sequences of independent and almost identically distributed variables with zero mean. We denote the corresponding increments as

$$r_i(t) = \epsilon_i(t) \quad \forall i, \quad (84)$$

where the  $\epsilon_i$  are random variables. In a traditional factor analysis, the above scenario takes the form of a ‘zero-factor’ model. Under this model, the aggregate increment over  $N$  stocks is expected to be narrowly distributed around

$$\{r_i(t)\} = \frac{1}{N} \sum_{i=1}^N \epsilon_i(t) \approx 0. \quad (85)$$

When  $\{r_i(t)\}$  takes small values around zero, we know from figure 12 that  $\{x_i(t)\}$  also takes small values around zero. Indeed, shuffled time series are in the linear regime that spans the range where the binary increment  $\{x_i(t)\}$  is consistent with a uniform random walk (see figure 9). Therefore we find that the zero-factor model (for the non-binary returns) and the uniform random walk (for the binary returns) are consistent with each other in the linear regime. In other words, when in our analysis we measure a value of  $\{x_i(t)\}$  that is consistent with a uniform random walk, we know that the original log-returns are consistent with a zero-factor model.

Next, we consider a one-factor model, where there is one dominant underlying factor assumed to control the dynamics of all the time series. In such a case, each return can be decomposed as

$$r_i(t) = \alpha_i \Phi_0(t) + \epsilon_i(t) \quad \forall i, \quad (86)$$

where  $\alpha_i$  is the ‘factor loading’ of the  $i$ th time series with the dominant factor  $\Phi_0(t)$ . When referring to stocks, the factor  $\Phi_0(t)$  is attributed to the market mode. It is known that, during crisis times when the markets are highly correlated, a one-factor model can describe the dynamics quite well. Under this model, the aggregate increment is

$$\{r_i(t)\} = \frac{1}{N} \sum_{i=1}^N \alpha_i \Phi_0(t) + \frac{1}{N} \sum_{i=1}^N \epsilon_i(t) \approx \{\alpha_i\} \Phi_0(t), \quad (87)$$

where  $\{\alpha_i\} \equiv \frac{1}{N} \sum_{i=1}^N \alpha_i$  is the average loading, which is independent of both  $i$  and  $t$ . This result implies that, when the market is well described by a one-factor model, the average increment  $\{r_i(t)\}$  that we measure in our analysis is proportional to the factor  $\Phi_0(t)$  itself. We note that the one-factor model is somehow similar to our biased random walk model, as it assumes a common drive for all the stocks. However, since  $\Phi_0(t)$  is fitted on the data, the one-factor model cannot distinguish between an endogenous or exogenous nature of the common drive. This situation is similar to when we use the observed value of  $\{r_i(t)\}$  as the driving field of the biased random walk (see section 6.2).

In financial analysis, the factor model can be used to filter the original time series and remove the one-factor component from them. When the model is a good approximation to the real market, the filtered returns are  $r_i(t) \approx \epsilon_i(t)$ , leading us back to equation (84) and the related considerations. In such a scenario, there is no correlation among the stocks, and each stock is acting as an i.i.d. variable. We therefore expect that, if we remove the market mode from the original time series, then (in periods where the market is indeed dominated by a single factor) we would obtain results similar to the shuffled case, and we would find the system in the uncoordinated phase of figure 9.

However, despite the fact that in certain conditions the one-factor model can generate the market behaviour, the model is too simplistic [3]. In reality the dynamics is more complex and can be attributed to many factors, that sometimes overlap with industrial (sub)sectors. Generally the different factors are identified by the largest, non-random eigenvalues of the empirical cross-correlation matrix, where the market mode relates to the highest eigenvalue [3]. The presence of many deviating eigenvalues is an indication of the fact that the one-factor model should be rejected. A more realistic,  $M$ -factor model is

$$r_i(t) = \sum_{j=0}^M \alpha_{ij} \Phi_j(t) + \epsilon_i(t) \quad \forall i, \quad (88)$$

where  $j = 0$  denotes a common market-wide factor as above, while  $j > 0$  denotes sector-specific factors. In such a case, our measured value of  $\{r_i(t)\}$  is

$$\{r_i(t)\} = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M \alpha_{ij} \Phi_j(t) + \frac{1}{N} \sum_{i=1}^N \epsilon_i(t) \approx \sum_{j=1}^M \{\alpha_{ij}\} \Phi_j(t) \quad (89)$$

which is a linear combination of the multiple factors controlling the market dynamics.



It should be noted that factor models cannot distinguish between an endogenous and exogenous origin for the factors  $\Phi_j(t)$  themselves, even if we invoke some information-theoretic criterion to rank different specifications of these models. By contrast, our binary models allow us to discriminate among these multiple scenarios, as we have shown in figure 9 and related discussion. Moreover, while our approach allows us to relate binary and non-binary increments of real time series and replicate the observed relationships among them (see figures 12 and 13), factor models cannot lead to a similar result, because they do not allow for a binary description.

## 7. Conclusions

We presented a novel method for the analysis of single and multiple binary time series. Our information-theoretic approach allowed us to extract and quantify the amount of information encoded in simple, empirically measured properties. This resulted in the possibility to associate an entropy value to a time series given its measured properties, and to compare the informativeness of different measured properties.

By employing our formalism, we have identified distinct regimes in the collective behavior of groups of stocks, corresponding to different levels of coordination that only depend on the average return of the binary time series. In each regime the market exhibits a dominant character: the market mode can be interpreted as an exogenous factor, as pure noise, or as a combination of endogenous and exogenous components. Moreover, each regime is characterized by the most informative property.

Finally and more importantly, we were able to replicate the observed nonlinear relations between binary and non-binary aggregate increments of real multiple time series. We have mathematically characterized these relations accurately, and interpreted them as the result of the fact that very large log-returns occur more often when most stocks are synchronized, i.e. when their increments have a common sign. Our findings suggest that the binary signatures carry significant information, and even allow one to measure the level of coordination in a way that is unaccessible to standard non-binary analyses.

## Acknowledgments

We thank Marc van Kralingen for a thorough reading of our manuscript and for identifying some mistakes. We acknowledge support from the Dutch Econophysics Foundation (Stichting Econophysics, Leiden, The Netherlands) with funds from beneficiaries of Duyfken Trading Knowledge BV, Amsterdam, The Netherlands. This work was also supported by the EU project MULTIPLEX (contract 317532) and The Netherlands Organization for Scientific Research (NWO/OCW).

## Appendix A. Models for single time series

We consider the case  $N = 1$ , i.e. when  $\mathbf{X}$  is a  $1 \times T$  matrix or equivalently a  $T$ -dimensional row vector. Let us denote the entries of  $\mathbf{X}$  as  $x(t)$ .



### A.1. Uniform random walk model

The trivial model is obtained when no constraints are enforced. In this case, there is no free parameter and the Hamiltonian has the form

$$H(\mathbf{X}) = 0. \quad (\text{A.1})$$

As a result, the partition function is

$$Z = \sum_{\mathbf{X}} 1 = 2^T \quad (\text{A.2})$$

which is nothing but the number of possible binary time series of length  $T$ . The probability of occurrence of a time series  $\mathbf{X}$  is then

$$P(\mathbf{X}) = \frac{1}{Z} = 2^{-T} \quad (\text{A.3})$$

and is completely uniform over the ensemble of all binary time series of length  $T$ . All the  $T$  elements of  $\mathbf{X}$  are mutually independent and identically distributed with probability

$$P_t(x) \equiv \text{Prob}(x(t) = x) = \begin{cases} 1/2 & x = -1 \\ 1/2 & x = +1 \end{cases} \quad (\text{A.4})$$

This results in a completely uniform random walk with zero expected value for each increment:

$$\langle x(t) \rangle = 0. \quad (\text{A.5})$$

While the (ensemble) variance of each increment equals

$$\text{Var}[x(t)] \equiv \langle x^2(t) \rangle - \langle x(t) \rangle^2 = 1. \quad (\text{A.6})$$

### A.2. Biased random walk model

We now consider the total increment as the simplest non-trivial (one-dimensional) constraint:

$$C(\mathbf{X}) = T \cdot M_1(\mathbf{X}) = T \cdot \overline{x(t)}. \quad (\text{A.7})$$

If we denote the corresponding (scalar) Lagrange multiplier by  $\theta$ , the Hamiltonian has the form

$$H(\mathbf{X}, \theta) = \theta \cdot T \cdot \overline{x(t)} = \theta \sum_{t=1}^T x(t). \quad (\text{A.8})$$

The partition function is

$$\begin{aligned} Z(\theta) &= \sum_{\mathbf{X}} e^{-\theta \sum_{t=1}^T x(t)} = \sum_{\mathbf{X}} \prod_{t=1}^T e^{-\theta x(t)} \\ &= \prod_{t=1}^T \sum_{x=\pm 1} e^{-\theta x} = \prod_{t=1}^T [e^{-\theta} + e^{+\theta}], \\ &= [e^{-\theta} + e^{+\theta}]^T \end{aligned} \quad (\text{A.9})$$

where, when interchanging the order of the sum and product, we have replaced the sum over all time series  $\mathbf{X}$  with the sum over the two possible values  $x = \pm 1$  of each individual entry.

The probability of the occurrence of a time series  $\mathbf{X}$  is

$$\begin{aligned} P(\mathbf{X} | \theta) &= \frac{e^{-\theta \sum_{t=1}^T x(t)}}{\left[ e^{-\theta} + e^{+\theta} \right]^T} = \prod_{t=1}^T \frac{e^{-\theta x(t)}}{e^{-\theta} + e^{+\theta}}, \\ &= \prod_{t=1}^T P_t(x(t) | \theta) \end{aligned} \quad (\text{A.10})$$

where we have introduced the probability  $P_t(x | \theta)$  of a given increment  $x = \pm 1$  at time  $t$ , which we identify as

$$P_t(x | \theta) = \frac{e^{-\theta x}}{e^{-\theta} + e^{+\theta}}. \quad (\text{A.11})$$

The above expression shows that the stochastic process corresponding to this model is a biased random walk, as the two outcomes  $x = \pm 1$  have a different probability, unless  $\theta = 0$  (which leads us back to the uniform random walk model considered above).

The expected value of the  $t$ th increment  $x(t)$  (representing the bias of the random walk) is

$$\langle x(t) \rangle_\theta = \sum_{x=\pm 1} x P_t(x | \theta) = \frac{e^{-\theta} - e^{+\theta}}{e^{-\theta} + e^{+\theta}} = -\tanh \theta \quad (\text{A.12})$$

and the variance is

$$\text{Var}[x(t)] = \langle x^2(t) \rangle_\theta - \langle x(t) \rangle_\theta^2 = 1 - \tanh^2 \theta. \quad (\text{A.13})$$

The maximum likelihood condition (16), fixing the value  $\theta^*$  of the parameter  $\theta$  given a real time series  $\mathbf{X}^*$ , reads

$$T \langle \overline{x(t)} \rangle = \sum_{t=1}^T \langle x(t) \rangle = -T \tanh \theta = T \cdot \overline{x^*(t)}, \quad (\text{A.14})$$

where  $\overline{x^*(t)}$  is the measured average increment in the observed time series  $\mathbf{X}^*$ . This yields

$$-\tanh \theta^* = \overline{x^*(t)} \quad (\text{A.15})$$

which gives a parameter value

$$\theta^* = -\text{artanh} \left[ \overline{x^*(t)} \right] = -\frac{1}{2} \ln \left[ \frac{1 + \overline{x^*(t)}}{1 - \overline{x^*(t)}} \right]. \quad (\text{A.16})$$

### A.3. One-dimensional Ising model

We now consider a model where, besides the constraint on the total increment specified in equation (33), we enforce an additional constraint on the time-delayed (lagged) quantity  $T \cdot B_1(\mathbf{X})$ , where  $B_1(\mathbf{X})$  is defined in equation (27) with  $\tau = 1$ . This amounts to enforcing the average one-step temporal autocorrelation of the time series. The resulting two-dimensional constraint can be written as the column vector

$$\vec{C}(\mathbf{X}) = \begin{pmatrix} C_1(\mathbf{X}) \\ C_2(\mathbf{X}) \end{pmatrix} = T \cdot \begin{pmatrix} M_1(\mathbf{X}) \\ B_1(\mathbf{X}) \end{pmatrix}. \quad (\text{A.17})$$

If we write the corresponding Lagrange multiplier as

$$\vec{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = - \begin{pmatrix} I \\ K \end{pmatrix}, \quad (\text{A.18})$$

then the Hamiltonian reads

$$\begin{aligned} H(\mathbf{X}, I, K) &= \vec{\theta} \cdot \vec{C}(\mathbf{X}) = T\theta_1 M_1(\mathbf{X}) + T\theta_2 B_1(\mathbf{X}) \\ &= -I \sum_{t=1}^T x(t) - K \sum_{t=1}^T x(t)x(t+1), \end{aligned} \quad (\text{A.19})$$

where we consider a periodicity condition as in equation (28) with  $\tau = 1$ , i.e.  $x(T+1) \equiv x(1)$ . Note that, when  $\mathbf{X}$  is a real binary time series of length  $T$ , this condition can be always enforced by adding one last (fictious) timestep  $T+1$  and a corresponding increment  $x(T+1)$  chosen equal to  $x(1)$ . For long time series, this has a negligible effect.

The above Hamiltonian coincides with that for the one-dimensional Ising model with periodic boundary conditions [46]. Each time step  $t$  is seen as a site in an ordered chain of length  $T$ , and each value  $x(t) = \pm 1$  is seen as the value of a spin sitting at that site. The model is analytically solvable, which allows us to apply it to real time series in our formalism. For the readers familiar with time series analysis but not necessarily with the Ising model, we briefly recall the standard solution of the model, adapting it from [46].

Applying the periodicity condition of equation (28) ensures that all sites (time steps) are statistically equivalent, i.e.:

$$\langle x(1) \rangle = \langle x(2) \rangle = \dots = \langle x(T) \rangle \quad (\text{A.20})$$

so that the system is translationally (here, temporally) invariant. The partition function is

$$Z(I, K) = \sum_{\mathbf{X}} \exp \left[ I \sum_{t=1}^T x(t) + K \sum_{t=1}^T x(t)x(t+1) \right]$$

and can be rewritten as a product of terms involving only two successive time steps:

$$Z(I, K) = \sum_{\mathbf{X}} \prod_{t=1}^T V(x(t), x(t+1)), \quad (\text{A.21})$$

where we have introduced the function  $V(x, y)$  defined as

$$V(x, y) \equiv \exp \left( I \frac{x+y}{2} + Kxy \right). \quad (\text{A.22})$$

We since both  $x$  and  $y$  can take only the values  $\pm 1$ , we can regard  $V(x, y)$  as the element of a  $2 \times 2$  matrix  $\mathbf{V}$  called the *transfer matrix* [46]:

$$\mathbf{V} \equiv \begin{pmatrix} V(+1, +1) & V(+1, -1) \\ V(-1, +1) & V(-1, -1) \end{pmatrix} = \begin{pmatrix} e^{K+I} & e^{-K} \\ e^{-K} & e^{K-I} \end{pmatrix}. \quad (\text{A.23})$$

This allows us to rewrite equation (A.21) as

$$Z(I, K) = \text{Tr}(\mathbf{V}^T). \quad (\text{A.24})$$

Let  $\vec{v}_1, \vec{v}_2$  denote the two eigenvectors of  $\mathbf{V}$ , and  $\lambda_1, \lambda_2$  the corresponding eigenvalues, so that

$$\mathbf{V}\vec{v}_j = \lambda_j\vec{v}_j, \quad j = 1, 2. \quad (\text{A.25})$$

The  $2 \times 2$  matrix  $\mathbf{Q} \equiv (\vec{v}_1, \vec{v}_2)$  (having column vectors  $\vec{v}_1$  and  $\vec{v}_2$ ) diagonalizes  $\mathbf{V}$ , i.e.

$$\mathbf{V} = \mathbf{Q} \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \mathbf{Q}^{-1}, \quad (\text{A.26})$$

where a direct calculation of the eigenvalues and eigenvectors yields

$$\lambda_1 = e^K \cosh I + \sqrt{e^{2K} \sinh^2 I + e^{-2K}} \quad (\text{A.27})$$

$$\lambda_2 = e^K \cosh I - \sqrt{e^{2K} \sinh^2 I + e^{-2K}} \quad (\text{A.28})$$

and

$$\mathbf{Q} = \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix}, \quad (\text{A.29})$$

with  $\phi$  defined by

$$\cot 2\phi \equiv e^{2K} \sinh I. \quad (\text{A.30})$$

It then follows that equation (A.24) simply reduces to

$$Z(I, K) = \text{Tr} \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}^T = \lambda_1^T + \lambda_2^T, \quad (\text{A.31})$$

and the probability of occurrence of a time series  $\mathbf{X}$  is

$$P(\mathbf{X} | I, K) = \frac{\prod_{t=1}^T V(x(t), x(t+1))}{\lambda_1^T + \lambda_2^T}. \quad (\text{A.32})$$

The above results allow us to analytically obtain expected values. That of  $x(t)$  is

$$\langle x(t) \rangle = \sum_{\mathbf{X}} x(t) P(\mathbf{X} | I, K) = \frac{\text{Tr}(\mathbf{S}\mathbf{V}^T)}{\lambda_1^T + \lambda_2^T}, \quad (\text{A.33})$$

where we have introduced the diagonal matrix

$$\mathbf{S} \equiv \begin{pmatrix} S(+1, +1) & S(+1, -1) \\ S(-1, +1) & S(-1, -1) \end{pmatrix} = \begin{pmatrix} +1 & 0 \\ 0 & -1 \end{pmatrix} \quad (\text{A.34})$$

having elements

$$S(x, y) \equiv x\delta(x, y). \quad (\text{A.35})$$

Similarly, for  $0 < s - t < T$  the expected value of  $x(t)x(s)$  is

$$\begin{aligned}\langle x(t)x(s) \rangle &= \sum_{\mathbf{X}} x(t)x(s)P(\mathbf{X} | I, K) \\ &= \frac{\text{Tr}(\mathbf{S}\mathbf{V}^{s-t}\mathbf{S}\mathbf{V}^{T+t-s})}{\lambda_1^T + \lambda_2^T}.\end{aligned}\quad (\text{A.36})$$

In the limit  $T \rightarrow \infty$  (corresponding to long time series in our case) with  $s - t$  fixed, these expressions become

$$\langle x(t) \rangle = \cos 2\phi \quad (\text{A.37})$$

$$\langle x(t)x(s) \rangle = \cos^2 2\phi + \sin^2 2\phi \left( \frac{\lambda_1}{\lambda_2} \right)^{s-t}. \quad (\text{A.38})$$

Now, we note that equations (A.33) and (A.36) manifestly show the translational (temporal) invariance of the model, as  $\langle x(t) \rangle$  is independent of  $t$  and  $\langle x(t)x(s) \rangle$  depends on  $t$  and  $s$  only through their difference  $s - t$ . This implies that, writing  $\tau \equiv s - t$  and performing a temporal average,

$$\langle M_1 \rangle = \cos 2\phi \quad (\text{A.39})$$

$$\langle B_\tau \rangle = \cos^2 2\phi + \sin^2 2\phi \left( \frac{\lambda_1}{\lambda_2} \right)^\tau. \quad (\text{A.40})$$

Using equation (A.30) we can rewrite these expressions in terms of the model parameters,  $I$  and  $K$ , as

$$\langle M_1 \rangle = \frac{e^{2K} \sinh I}{\sqrt{1 + e^{4K} \sinh^2 I}} \quad (\text{A.41})$$

$$\langle B_\tau \rangle = \frac{e^{4K} \sinh^2 I + (\lambda_1/\lambda_2)^\tau}{1 + e^{4K} \sinh^2 I}. \quad (\text{A.42})$$

The expected value of the autocorrelation defined in equation (47) can be approximated as the ratio of two expected values as follows:

$$\langle A_\tau \rangle \equiv \left\langle \frac{B_\tau - M_1^2}{1 - M_1^2} \right\rangle \approx \frac{\langle B_\tau \rangle - \langle M_1^2 \rangle}{1 - \langle M_1^2 \rangle} = \left( \frac{\lambda_1}{\lambda_2} \right)^\tau. \quad (\text{A.43})$$

## Appendix B. Models for single cross-sections of multiple time series

For a single cross-section of a set of  $N$  multiple time series,  $\mathbf{X}$  is a  $N \times 1$  matrix or equivalently a  $N$ -dimensional column vector. We denote the entries of  $\mathbf{X}$  as  $x_i$ .

### B.1. Uniform random walk model

The uniform random walk is a simple modification of the same model that we considered for single time series, where  $x(t)$  is replaced by  $x_i$  and  $T$  is replaced by  $N$ . This model is obtained when no constraints are enforced. The Hamiltonian is

$$H(\mathbf{X}) = 0 \quad (\text{B.1})$$

and the partition function is simply the number of possible configurations for a single cross-section of  $N$  stocks:

$$Z = \sum_{\mathbf{X}} 1 = 2^N. \quad (\text{B.2})$$

The probability of occurrence of a cross-section  $\mathbf{X}$  is

$$P(\mathbf{X}) = \frac{1}{Z} = 2^{-N} \quad (\text{B.3})$$

and is completely uniform over the ensemble of all cross-sections of  $N$  stocks. All the  $N$  elements of  $\mathbf{X}$  are mutually independent and identically distributed with probability

$$P_i(x) \equiv \text{Prob}(x_i = x) = \begin{cases} 1/2 & x = -1 \\ 1/2 & x = +1 \end{cases}. \quad (\text{B.4})$$

This results in a completely uniform random walk with zero expected value

$$\langle x_i \rangle = 0 \quad (\text{B.5})$$

and maximum variance

$$\text{Var}[x_i] \equiv \langle x_i^2 \rangle - \langle x_i \rangle^2 = 1. \quad (\text{B.6})$$

### B.2. Biased random walk model

Also this model is analogous to the corresponding model for single time series. We select the total daily increment of the cross-section  $\mathbf{X}$  as the constraint:

$$C(\mathbf{X}) = N \cdot M_1(\mathbf{X}) = N \cdot \{x_i\}. \quad (\text{B.7})$$

Let the corresponding Lagrange multiplier be denoted by  $\theta$ . The Hamiltonian is

$$H(\mathbf{X}, \theta) = \theta \cdot N \cdot \{x_i\} = \theta \sum_{i=1}^N x_i \quad (\text{B.8})$$

and the partition function is

$$\begin{aligned} Z(\theta) &= \sum_{\mathbf{X}} e^{-\theta \sum_{i=1}^N x_i} = \sum_{\mathbf{X}} \prod_{i=1}^N e^{-\theta x_i} \\ &= \prod_{i=1}^N \sum_{x=\pm 1} e^{-\theta x} = \prod_{i=1}^N [e^{-\theta} + e^{+\theta}] \\ &= [e^{-\theta} + e^{+\theta}]^N. \end{aligned} \quad (\text{B.9})$$

The probability of the occurrence of a cross-section  $\mathbf{X}$  is

$$\begin{aligned} P(\mathbf{X} | \theta) &= \frac{e^{-\theta \sum_{i=1}^N x_i}}{\left[ e^{-\theta} + e^{+\theta} \right]^N} = \prod_{i=1}^N \frac{e^{-\theta x_i}}{e^{-\theta} + e^{+\theta}}, \\ &= \prod_{i=1}^N P_i(x_i | \theta) \end{aligned} \quad (\text{B.10})$$

where we have introduced the probability  $P_i(x | \theta)$  of a given increment  $x = \pm 1$  for stock  $i$ , which we identify as

$$P_i(x | \theta) = \frac{e^{-\theta x}}{e^{-\theta} + e^{+\theta}}. \quad (\text{B.11})$$

Just like the corresponding model for single time series, this model is a biased random walk, because the two outcomes  $x = \pm 1$  have a different probability unless  $\theta = 0$ .

The expected value of the  $i$ th increment  $x_i$  is

$$\langle x_i \rangle_\theta = \sum_{x=\pm 1} x P_i(x | \theta) = \frac{e^{-\theta} - e^{+\theta}}{e^{-\theta} + e^{+\theta}} = -\tanh \theta \quad (\text{B.12})$$

and the variance is

$$\text{Var}[x_i] = \langle x_i^2 \rangle_\theta - \langle x_i \rangle_\theta^2 = 1 - \tanh^2 \theta. \quad (\text{B.13})$$

The maximum likelihood condition (16), fixing the value  $\theta^*$  of the parameter  $\theta$  given a real cross-section  $\mathbf{X}^*$ , reads

$$N \langle \{x_i\} \rangle = \sum_{i=1}^N \langle x_i \rangle = -N \tanh \theta = N \cdot \{x_i^*\}, \quad (\text{B.14})$$

where  $\{x_i^*\}$  is the measured average increment of the observed cross-section  $\mathbf{X}^*$ . This yields

$$-\tanh \theta^* = \{x_i^*\} \quad (\text{B.15})$$

which gives a parameter value

$$\theta^* = -\text{artanh} \left[ \{x_i^*\} \right] = -\frac{1}{2} \ln \left[ \frac{1 + \{x_i^*\}}{1 - \{x_i^*\}} \right]. \quad (\text{B.16})$$

### B.3. Mean-field Ising model

In this model, we enforce two constraints: the total increment and the total coupling between stocks. The resulting two-dimensional constraint can be written as

$$\vec{C}(\mathbf{X}) = \begin{pmatrix} C_1(\mathbf{X}) \\ C_2(\mathbf{X}) \end{pmatrix} = \begin{pmatrix} N \cdot M_1(\mathbf{X}) \\ D(\mathbf{X}) \end{pmatrix}. \quad (\text{B.17})$$



We can write the corresponding Lagrange multiplier as

$$\vec{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = - \begin{pmatrix} h \\ J \end{pmatrix} \quad (\text{B.18})$$

and the Hamiltonian as

$$H(\mathbf{X}, h, J) = -h \sum_{i=1}^N x_i - J \sum_{i<j} x_i x_j. \quad (\text{B.19})$$

Note that here we are not enforcing nearest-neighbor interactions as in the one-lagged model for single time series, but market-wide interactions among all stocks for the same time step (cross-section). This is the result of the fact that, when considering cross-sections, there is no natural notion of ‘lattice sites’ induced by e.g. a temporal ordering as in the one-lagged model. In other words, pairs of stocks in a cross-section are neither ‘close’ nor ‘distant’. We therefore assume a common interaction strength  $J$  among all stocks.

The above model, known as the mean-field Ising model, is analytically solvable. Here we adapt the derivation illustrated in [46]. We first note that, since  $x_i^2 = 1$  for all  $i$ ,  $H(\mathbf{X}, h, J)$  can be expressed as a function of  $M_1(\mathbf{X})$  alone:

$$H(\mathbf{X}, h, J) = -hNM_1(\mathbf{X}) - \frac{J}{2} [N^2 M_1^2(\mathbf{X}) - N]. \quad (\text{B.20})$$

This implies that the sum over configurations in the partition function can be replaced by a sum over the allowed values of  $M_1(\mathbf{X})$ , weighted by the number of configurations for each value. If we denote by  $r$  the number of increments that are negative ( $x = -1$ ), and by  $(N - r)$  the number of increments that are positive ( $x = +1$ ), then we can write the Hamiltonian as a function of  $r$  alone through the expression

$$NM_1(\mathbf{X}) = N - 2r. \quad (\text{B.21})$$

The partition function can therefore be calculated as

$$Z(h, J) \equiv \sum_{\mathbf{X}} e^{-H(\mathbf{X}, h, J)} = \sum_{r=1}^N C_r, \quad (\text{B.22})$$

where

$$C_r \equiv \frac{N!}{r!(N-r)!} e^{h(N-2r) + \frac{J}{2}[N(N-2r)^2 - N]} \quad (\text{B.23})$$

incorporates the binomial coefficient enumerating the configurations with given  $r$ . The expected increment is therefore

$$\langle M_1 \rangle = \left\langle 1 - \frac{2r}{N} \right\rangle = \frac{\sum_{r=1}^N \left(1 - \frac{2r}{N}\right) C_r}{Z(h, J)} \quad \forall i. \quad (\text{B.24})$$

When  $N$  is large, a traditional derivation [46] shows that the sum at the numerator of equation (B.24) is dominated by the single addendum corresponding to the maximum of  $C_r$ . The same applies to the partition function at the denominator. If  $r_0$  denotes the value of  $r$  such that  $C_r$  is maximum, we then get

$$\langle M_1 \rangle \approx 1 - \frac{2r_0}{N}. \quad (\text{B.25})$$

A further expansion [46] finally shows that, given  $h$  and  $J$ , the expected value  $\langle M_1 \rangle$  is the solution of the nonlinear equation

$$\langle M_1 \rangle = \tanh[(N-1)J\langle M_1 \rangle + h]. \quad (\text{B.26})$$

From the above equation, one can infer the existence of a phase transition in the model, separating a regime where the expected ‘magnetization’ (here the average increment  $\langle M_1 \rangle$ ) is zero from one where it is non-zero [46]. This transition is discussed in section 5.3.

Before proceeding further, we note a peculiarity of the model, which has implications for the applicability of our maximum likelihood approach. An argument similar to that leading to equation (B.25) implies that the second moment of  $M_1(\mathbf{X})$  can be expressed as

$$\langle M_1^2 \rangle = \left\langle \left(1 - \frac{2r}{N}\right)^2 \right\rangle \approx \left(1 - \frac{2r_0}{N}\right)^2 \approx \langle M_1 \rangle^2. \quad (\text{B.27})$$

This implies that

$$\text{Var}[M_1] \equiv \langle M_1^2 \rangle - \langle M_1 \rangle^2 = 0, \quad (\text{B.28})$$

or in other words that  $M_1(\mathbf{X})$  is no longer a random variable. As a consequence, something unusual happens when we apply the maximum likelihood principle. From equation (B.20), and recalling the general result embodied by equation (20) in section 3.2, it is clear that the parameter values  $h^*$  and  $J^*$  maximizing the likelihood can be found as the solution to the two coupled equations

$$\langle M_1 \rangle = M_1(\mathbf{X}^*) \quad (\text{B.29})$$

$$\langle M_1^2 \rangle = M_1^2(\mathbf{X}^*). \quad (\text{B.30})$$

However, equation (B.27) implies that equation (B.30) can be rewritten as

$$\langle M_1 \rangle^2 = M_1^2(\mathbf{X}^*) \quad (\text{B.31})$$

which coincides with equation (B.29). So equations (B.29) and (B.30) are equivalent, and they cannot be used to uniquely determine the two unknown parameters  $h^*$  and  $J^*$ . This is the result of the fact that, when fitted to the data, the model is actually over-constrained: there are two parameters to fit the only constraint ( $M_1$ ) on which the Hamiltonian depends. This aspect of the model is not manifest when  $M_1$  is regarded as a function of  $h$  and  $J$ , as usually done when simulating spin systems.

The above consideration implies that we should drop one of the two parameters and consider the two cases  $J = 0$  and  $h = 0$  separately. The former case coincides with the biased random walk model that we already discussed, and we will not discuss it any further. The latter case will instead represent our genuine specification of the ‘mean-field’ model. Setting  $h = 0$  implies

$$H(\mathbf{X}, 0, J) = -\frac{J}{2} \left[ N^2 M_1^2(\mathbf{X}) - N \right] \quad (\text{B.32})$$

and

$$\langle M_1 \rangle = \tanh \left[ (N-1)J \langle M_1 \rangle \right]. \quad (\text{B.33})$$

Applying the maximum likelihood principle to equation (B.32) tells us to select  $J^*$  as the solution of equation (B.30). However, we have seen that this condition leads to equation (B.31), which is actually equivalent to equation (B.29). Therefore, the value of  $J^*$  can be found by replacing  $\langle M_1 \rangle$  with the observed value  $M_1(\mathbf{X}^*) = \{x_i^*\}$  in equation (B.33), which leads to

$$\{x_i^*\} = \tanh \left[ (N-1)J^* \{x_i^*\} \right]. \quad (\text{B.34})$$

Note that in the traditional situation one is interested in finding the (expected) magnetization given a value of  $J$ , which implies that the transcendental equation (B.33) should be solved numerically. Here, we are instead facing the inverse situation where we look for the value of  $J^*$  given the (observed) value of the magnetization. In this quite unusual case, it turns out that equation (B.34) can be inverted to give the following analytical solution:

$$J^* = \frac{\text{artanh} \{x_i^*\}}{\{x_i^*\}(N-1)} = \frac{1}{2\{x_i^*\}(N-1)} \ln \left[ \frac{1 + \{x_i^*\}}{1 - \{x_i^*\}} \right]. \quad (\text{B.35})$$

Once this value is calculated, it can be inserted into the probability

$$P(\mathbf{X}^* | 0, J) = \frac{e^{-H(\mathbf{X}^*, 0, J)}}{Z(0, J)} = \frac{e^{JN(N\{x_i^*\}^2 - 1)/2}}{\sum_{r=1}^N \frac{N!}{r!(N-r)!} e^{J[(N-2r)^2 - N]/2}}, \quad (\text{B.36})$$

(where we have set  $h = 0$ ) to obtain the maximized likelihood of generating the observed cross-section  $\mathbf{X}^*$  under the mean-field model.

## References

- [1] Mantegna R N and Stanley H E 1999 *An Introduction to Econophysics: Correlation and Complexity in Finance* (Cambridge, UK: Cambridge University Press)
- [2] Sinha S, Chatterjee A, Chakraborti A and Chakrabarti B K 2010 *Econophysics: An Introduction* (New York: Wiley)
- [3] Bouchaud J-P and Potters M 2000 *Theory of Financial Risk: From Statistical Physics to Risk Management* (Cambridge: Cambridge University Press)
- [4] Bouchaud J P 2001 Power laws in economics and finance: some ideas from physics *Quant. Finance* **1** 105
- [5] Mantegna R N *et al* 1995 Scaling behaviour in the dynamics of an economic index *Nature* **376** 46–9  
Mantegna R N *et al* 1996 Turbulence and financial markets *Nature* **383** 587–8  
Plerou V *et al* 2003 Econophysics: two-phase behaviour of financial markets *Nature* **421** 130  
Gabaix X *et al* 2003 A theory of power-law distributions in financial market fluctuations *Nature* **423** 267–70
- [6] Schneidman E, Berry M J, Segev R and Bialek W 2006 Weak pairwise correlations imply strongly correlated network states in a neural population *Nature* **440** 1007–12
- [7] Dal'Maso Peron T and Rodrigues F 2011 Collective behavior in financial markets *Europhys. Lett.* **96** 48004
- [8] Mantegna R 1999 Hierarchical structure in financial markets *Eur. Phys. J. B* **11** 193–7

- [9] Onnela J, Chakraborti A, Kaski K and Kertesz J 2003 Dynamic asset trees and black Monday *Physica A* **324** 247–52
- [10] la Spada G, Farmer J and Lillo F 2008 The non-random walk of stock prices: the long-term correlation between signs and sizes *Eur. Phys. J. B* **64** 607–14
- [11] Petersen A M, Wang F, Havlin S and Stanley H E 2010 Stanley quantitative law describing market dynamics before and after interest rate change *Phys. Rev. E* **81** 066121
- [12] Black F 1976 Studies of stock price volatility changes *Proc. of the 1976 Meeting of the American Statistical Association, Business and Economics Statistics Section* (Washington, DC: American Statistical Association) pp 177–81
- [13] Andersen T G, Bollerslev T, Diebold F X and Vega C 2003 Micro effects of macro announcements: real-time price discovery in foreign exchange *Am. Econ. Rev.* **93** 38
- [14] Bekaert G and Wu G 2000 Asymmetric volatility and risk in equity markets *Rev. Financial Stud.* **13** 1–42
- [15] Bouchaud J-P, Matacz A and Potters M 2001 Leverage effect in financial markets: the retarded volatility model *Phys. Rev. Lett.* **87** 228701
- [16] Boguna M and Serrano M A 2005 Generalized percolation in random directed networks *Phys. Rev. E* **72** 016106
- [17] Hertz J, Roudi Y and Tyrcha J 2011 Ising models for inferring network structure from spike data arXiv:1106.1752
- Quenet B and Horn D 2003 The dynamic neural filter: a binary model of spatiotemporal coding *Neural Comput.* **15** 309–29
- [18] Jaynes E T 1957 Information theory and statistical mechanics *Phys. Rev.* **106** 620
- [19] Jaynes E T 1957 Information theory and statistical mechanics II *Phys. Rev.* **108** 171–90
- [20] Shannon C E 1948 A mathematical theory of communication *Bell Syst. Tech. J.* **27** 379–423 623–656
- [21] Borst A and Theunissen F E 1999 Information theory and neural coding *Nat. Neurosci.* **2** 947–57
- [22] Wasserman S 1994 *Social Network Analysis: Methods and Applications* (Cambridge: Cambridge University Press)
- [23] Park J and Newman M E J 2004 The statistical mechanics of networks *Phys. Rev. E* **70** 066117
- [24] Mehta M L 1960 On the statistical properties of the level-spacings in nuclear spectra *Nucl. Phys.* **18** 395–419
- Mehta M L and Dyson F J 1963 Statistical theory of the energy levels of complex systems. V *J. Math. Phys.* **4** 713
- Mehta M L 1971 A note on the correlations between eigenvalues of a random matrix *Commun. Math. Phys.* **20** 245–50
- [25] Mehta M L 1991 *Random Matrices* (Boston: Academic)
- [26] Wigner E P 1951 On a class of analytic functions from the quantum theory of collisions *Ann. Math.* **53** 36
- Wigner E P 1951 On the statistical distribution of the widths and spacings of nuclear resonance levels *Proc. Camb. Phil. Soc.* **47** 790–8
- [27] Dyson F J 1962 *J. Math. Phys.* **3** 140
- Dyson F J and Mehta M L 1963 *J. Math. Phys.* **4** 701
- [28] Sengupta A M and Mitra P P 1999 Distributions of singular values for some random matrices *Phys. Rev. E* **60** 3
- [29] Startz R 2008 Binomial autoregressive moving average models with an application to US recession *J. Business Economic Stat.* **26** 1–8
- [30] Kauppi H and Saikkonen P 2008 Predicting US recessions with dynamic binary response models *Rev. Economics Stat.* **90** 777–91
- [31] Preis T, Schneider J J and Stanley H E 2011 Switching processes in financial markets *Proc. Natl Acad. Sci. USA* **108** 7674–8
- [32] Squartini T and Garlaschelli D 2011 Analytical maximum-likelihood method to detect patterns in real networks *New J. Phys.* **13** 083001
- [33] Squartini T, Fagiolo G and Garlaschelli D 2011 Randomizing world trade. I. A binary network analysis *Phys. Rev. E* **84** 046117

- [34] Squartini T, Fagiolo G and Garlaschelli D 2011 Randomizing world trade. II. A weighted network analysis *Phys. Rev. E* **84** 046118
- [35] Akaike H 1974 A new look at the statistical model identification *IEEE Trans. Autom. Control* **19** 716–23
- [36] Bunde A and Havlin S 1994 *Fractals in Science* (Berlin: Springer)
- [37] Kolmogorov A N 1963 On tables of random numbers *Sankhyā: Indian J. Stat. Ser. A* **25** 369–76
- [38] Laloux L, Cizeau P, Bouchaud J-P and Potters M 1999 Noise dressing of financial correlation matrices *Phys. Rev. Lett.* **83** 1467
- [39] Plerou V, Gopikrishnan P, Rosenow B, Amaral L A N and Stanley H E 1999 Universal and nonuniversal properties of cross correlations in financial time series *Phys. Rev. Lett.* **83** 1471
- [40] Plerou V, Gopikrishnan P, Rosenow B, Amaral L A N and Stanley H E 2000 A random matrix theory approach to financial cross-correlations *Physica A* **287** 374
- [41] Plerou V, Gopikrishnan P, Rosenow B, Amaral L A N, Guhr T and Stanley H E 2002 Random matrix approach to cross-correlations in financial data *Phys. Rev. E* **64** 066126
- [42] Gopikrishnan P, Rosenow B, Plerou V and Stanley H E 2001 Quantifying and interpreting collective behavior in financial markets *Phys. Rev. E* **64** 035106
- [43] Rosenow B, Gopikrishnan P, Plerou V and Stanley H E 2002 Random magnets and correlation of stock price fluctuations *Physica A* **314** 762–7
- [44] Garlaschelli D and Loffredo M I 2008 Maximum likelihood: extracting unbiased information from complex networks *Phys. Rev. E* **78** 015101(R)
- [45] Burnham K P and Anderson D R 2002 *Model Selection and Multimodel Inference* 2nd edn (Berlin: Springer)
- [46] Baxter R J 2007 *Exactly Solved Models in Statistical Mechanics* (New York: Dover)
- [47] Sato A and Takayasu H 1998 Dynamic numerical models of stock market price: from microscopic determinism to macroscopic randomness *Physica A* **250** 231–52
- [48] Krawiecki A, Holyst J A and Helbing D 2002 Volatility clustering and scaling for financial time series due to attractor bubbling *Phys. Rev. Lett.* **89** 158701
- [49] Gontis V and Kaulakys B 2004 Multiplicative point process as a model of trading activity *Physica A* **343** 505–14