

A general statistical test for correlations in a finite-length time seriesJeffery A. Hanson and Haw Yang^{a)}Department of Chemistry, University of California at Berkeley, Berkeley, California 94720, USA and
Physical Biosciences Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA

(Received 27 March 2008; accepted 30 April 2008; published online 3 June 2008)

The statistical properties of the autocorrelation function from a time series composed of independently and identically distributed stochastic variables has been studied. Analytical expressions for the autocorrelation function's variance have been derived. It has been found that two common ways of calculating the autocorrelation, moving-average and Fourier transform, exhibit different uncertainty characteristics. For periodic time series, the Fourier transform method is preferred because it gives smaller uncertainties that are uniform through all time lags. Based on these analytical results, a statistically robust method has been proposed to test the existence of correlations in a time series. The statistical test is verified by computer simulations and an application to single-molecule fluorescence spectroscopy is discussed. © 2008 American Institute of Physics. [DOI: 10.1063/1.2931943]

I. INTRODUCTION

The dynamics of a system are usually characterized by the time correlation function of a variable q representing the system's state. If x is the experimental realization of q , then the time correlation function of interest is

$$C_{xx}(s) = \langle x(t)x(t+s) \rangle - \langle x \rangle^2. \quad (1)$$

where $\langle \dots \rangle$ represents ensemble averaging. This time correlation can be recovered from a prolonged, time-averaged study of a single system if the system is ergodic and when the duration of observation T approaches infinity

$$C_{xx}(s) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T x(t)x(t+\tau) dt - \left[\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T x(t) dt \right]^2. \quad (2)$$

In practice, the correlation function is calculated from an experimentally recorded time series, $[x_i]$, $i=1,\dots,N$, with each data point discretely sampled at a time interval δt .

There are two common ways of calculating the empirical correlation function. The moving-average approach is appropriate for time series that are potentially aperiodic and is given by

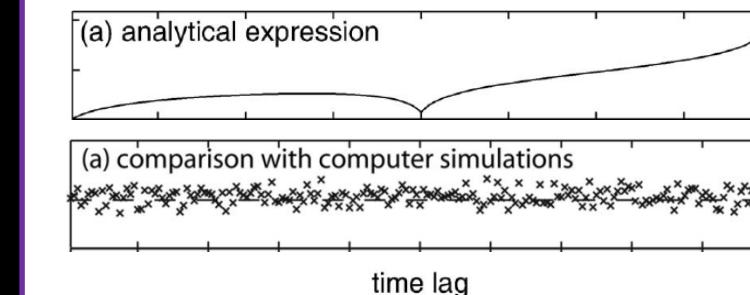
$$\bar{C}_{xx}(s) = \frac{1}{N-s} \sum_{i=1}^{N-s} x_i x_{i+s} - \bar{x}^2, \quad (3)$$

where $j_s = s/\delta t$ and $\bar{x} = \sum x_i/N$. If the time series can be assumed to exhibit a periodicity of N , i.e., $x_i = x_{i+N}$, then the correlation function can be calculated using discrete Fourier transform (DFT) and its inverse operation (IDFT). $\bar{C}_{xx}(s) = iDFT[\bar{x}_i^* \bar{x}_j]$, where $\bar{x}_i^* = DFT[x_i]$ denotes the i th Fourier-transform element, and $*$ indicates complex conjugate. The Fourier transform method is equivalent to calculating

^{a)}Electronic mail: hawyang@berkeley.edu.

A general statistical test for correlations in a finite-length time series

Jeffery A. Hanson and Haw Yang



time lag

Def. Time correation function

$$C_{xx}(s) = \langle x(t)x(t+s) \rangle - \langle x \rangle^2$$

Emperical corr. fn. w. discrete time series $\{x_i\}, i=1, \dots, N$

$$\bar{C}_{xx}(s) = \frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} - \bar{x}^2$$

The variance

$$\text{var}\{\bar{C}_{xx}(s)\} = \text{var}\left\{ \frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} \right\} + \text{var}\{\bar{x}^2\} - 2\text{cov}\left\{ \frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}, \bar{x}^2 \right\}$$

Covariance

$$\text{cov}(X, Y) = E[(X - E(X))(Y - E(Y))] = E[XY] - E[X]E[Y]$$

variance of the sum

$$\text{var}\{X + Y\} = \text{var}\{X\} + \text{var}\{Y\} - 2\text{cov}\{X, Y\}$$

The variance

$$\text{var}\{\bar{C}_{xx}(s)\} = \text{var}\left\{ \frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} \right\} + \text{var}\{\bar{x}^2\} - 2 \text{cov}\left\{ \frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}, \bar{x}^2 \right\}$$

$$\text{var}\left\{ \frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} \right\} = E\left\{ \left(\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} \right)^2 \right\} - E\left\{ \left(\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} \right) \right\}^2$$

for $1 < j_s < N/2$,

$$\begin{aligned} &= \frac{1}{N-j_s} E\{X_i^2 X_j^2\} + \frac{2(N-2j_s)}{(N-j_s)^2} E\{X_i X_j^2 X_k\} + \left[1 - \frac{3N-5j_s}{(N-j_s)^2} \right] E\{X_i X_j X_k X_l\} - E\{x_i\}^4 \\ &= \frac{1}{N-j_s} E\{X^2\}^2 + \frac{2(N-2j_s)}{(N-j_s)} E\{X^2\} E\{X\}^2 - \frac{3N-5j_s}{(N-j_s)^2} E\{X\}^4 \end{aligned}$$

for $N/2 < j_s < N$,

$$= \frac{1}{N-j_s} \left[E\{X^2\}^2 - E\{X\}^4 \right]$$

The variance

$$\text{var}\{\bar{C}_{xx}(s)\} = \text{var}\left\{\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}\right\} + \boxed{\text{var}\{\bar{x}^2\}} - \boxed{2 \text{cov}\left\{\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}, \bar{x}^2\right\}}$$

$$\begin{aligned} \boxed{\text{var}\{\bar{x}^2\}} &= \text{var}\left\{\left(\frac{1}{N} \sum_{i=1}^N x_i\right)^2\right\} = E\left\{\left(\frac{1}{N} \sum_{i=1}^N x_i\right)^4\right\} - E\left\{\left(\frac{1}{N} \sum_{i=1}^N x_i\right)^2\right\}^2 \\ &= \frac{1}{N^3} E\{X^4\} + \frac{4(N-1)}{N^3} E\{X^3\} E\{X\} + \frac{4(N-1)(N-3)}{N^3} E\{X^2\} \{X\}^2 - \frac{4N^2 - 10N + 6}{N^3} E\{X\}^4 \end{aligned}$$

$$\begin{aligned} \boxed{\text{cov}\left\{\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}, \bar{x}^2\right\}} \\ &= E\left\{\left[\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}\right] [\bar{x}^2]\right\} - E\left\{\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}\right\} E\{\bar{x}^2\} \\ &= \frac{2}{N^2} E\{X^3\} E\{X\} + \frac{2}{N^2} E\{X^2\}^2 + \frac{4N-10}{N^2} E\{X\}^2 E\{X^2\} - \frac{4N-6}{N^2} E\{X\}^4 \end{aligned}$$

The variance

$$\text{var}\{\bar{C}_{xx}(s)\} = \text{var}\left\{\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}\right\} + \text{var}\{\bar{x}^2\} - 2 \text{cov}\left\{\frac{1}{N-j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s}, \bar{x}^2\right\}$$

Result 1

$$\text{var}\{\bar{C}_{xx}(s)\} = h_1(N) + h_2(N, j_s),$$

$$\begin{aligned} h_1(N) &= \frac{1}{N^3} E\{X^4\} - \frac{4}{N^3} E\{X^3\} E\{X\} - \frac{2N+3}{N^3} E\{X^2\}^2 \\ &\quad - \frac{4(N^2-N-3)}{N^3} E\{X^2\} E\{X\}^2 \\ &\quad + \frac{2(2N-3)(N+1)}{N^3} E\{X\}^4, \end{aligned}$$

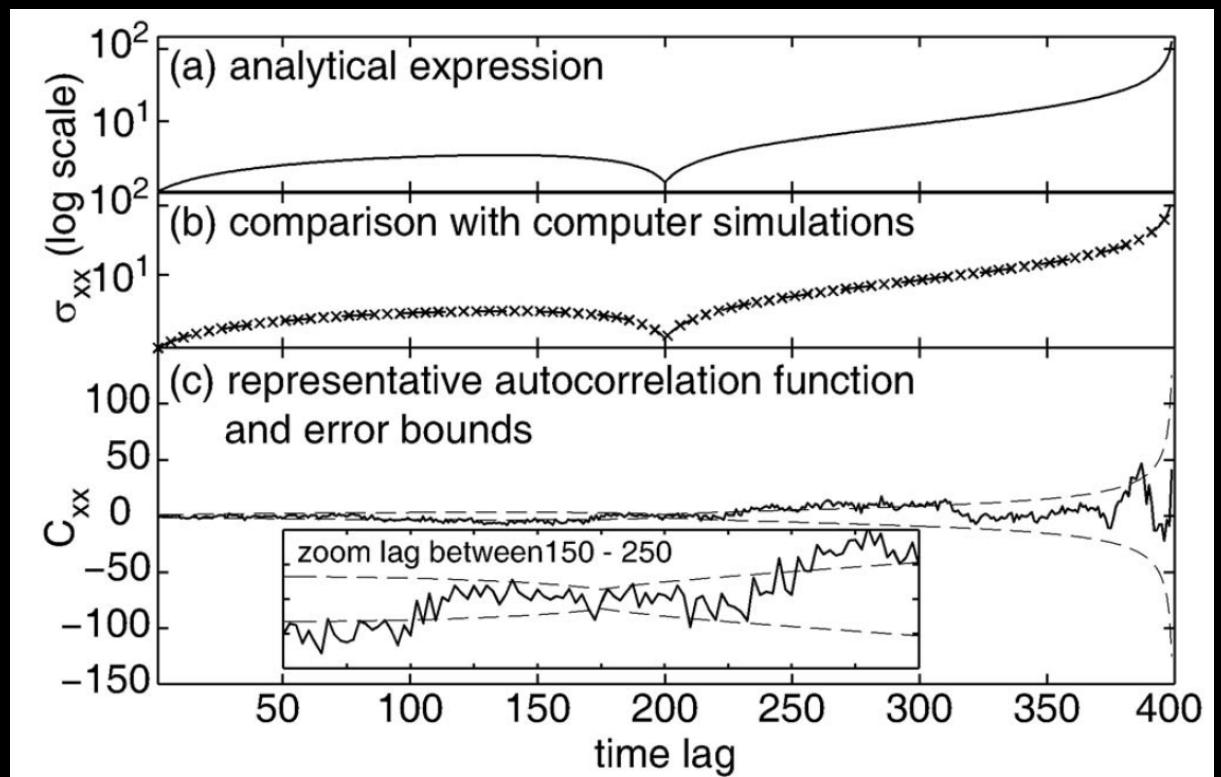
$$h_2(N, j_s) = \begin{cases} \frac{1}{N-j_s} E\{X^2\}^2 + \frac{2(N-2j_s)^2}{(N-j_s)^2} E\{X^2\} E\{X\}^2 - \frac{3N-5j_s}{(N-j_s)^2} E\{X^2\} E\{X\}^2, & 1 \leq j_s < (N/2); \\ \frac{1}{N-j_s} [E\{X^2\}^2 - E\{X\}^4], & (N/2) \leq j_s < N. \end{cases}$$

- (a) std from N=400
- (b) 10,000 simulated Poisson series (400-long)
- (c) an example of corr. fn. from simulation

For the Poisson series

$$h_1^{\text{poi}}(\lambda, N) = -\left(\frac{4\lambda^3}{N} + \frac{2\lambda^2}{N^2} - \frac{\lambda}{N^3}\right),$$

$$h_2^{\text{poi}}(\lambda, N, j_s) = \begin{cases} \frac{\lambda^2 + 6\lambda^3}{N - j_s} - \frac{2N\lambda^3}{(N - j_s)^2}, & 1 \leq j_s < (N/2); \\ \frac{\lambda^2 + 2\lambda^3}{N - j_s}, & (N/2) \leq j_s < N. \end{cases}$$

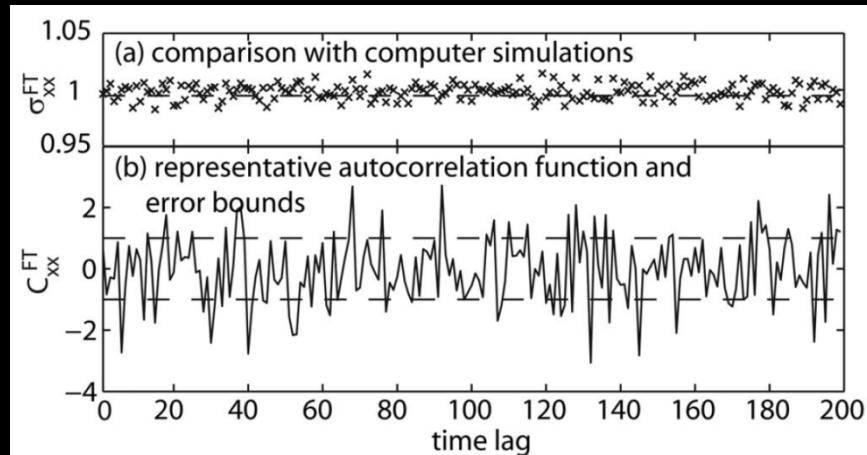


Fourier transform of the correation function

$$\bar{C}_{xx}(s) = iDFT \left\{ \tilde{x}_i^* \tilde{x}_j \right\}_s$$

$$\tilde{x}_j = DFT \left\{ x_i \right\}_j$$

$$\bar{C}_{xx}^{FT}(s) = \frac{1}{N} \sum_{i=1}^N x_i x_{i+j_s} - \bar{x}^2$$



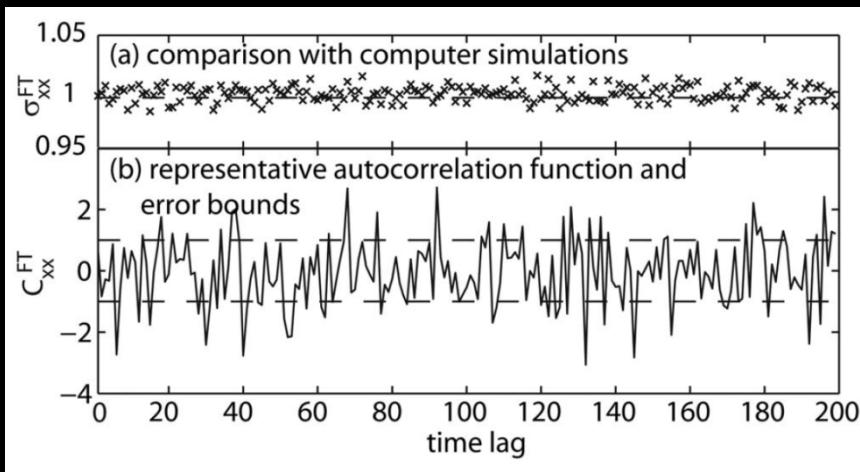
Result 2

$$\begin{aligned} \text{var}\{\bar{C}_{xx}^{FT}(s)\} &= \frac{1}{N^3} E\{X^4\} - \frac{4}{N^3} E\{X^3\} E\{X\} \\ &\quad + \frac{(N-3)(N+1)}{N^3} E\{X^2\}^2 \\ &\quad - \frac{2(N^2 - 2N - 6)}{N^3} E\{X^2\} E\{X\}^2 \\ &\quad + \frac{N^2 - 2N - 6}{N^3} E\{X\}^4. \end{aligned}$$

For the Poisson series

$$\text{var}\{\bar{C}_{xx}^{FT}(s)\}_{\text{poi}} = \frac{\lambda^2}{N} - \frac{2\lambda^2}{N^2} + \frac{\lambda}{N^3}$$

Statistical test for the existence of correlation



Result 3

$$Z_N = \left| \frac{1}{n_t} \sum_{j_s=1}^{n_t} \bar{C}_{xx}^{\text{FT}}(j_s) \right| > c_{1-\alpha},$$

$$c_{1-\alpha} = \sqrt{2} \operatorname{Erfc}^{-1}(\alpha) \bar{\sigma}_{xx}^{\text{FT}}$$

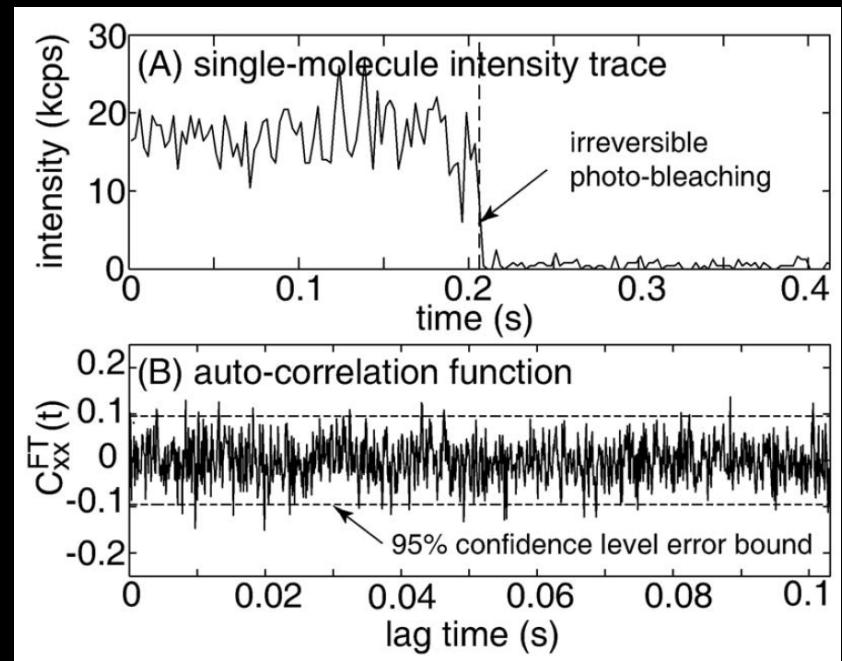
Statistical test for the existence of correlation

Result 3

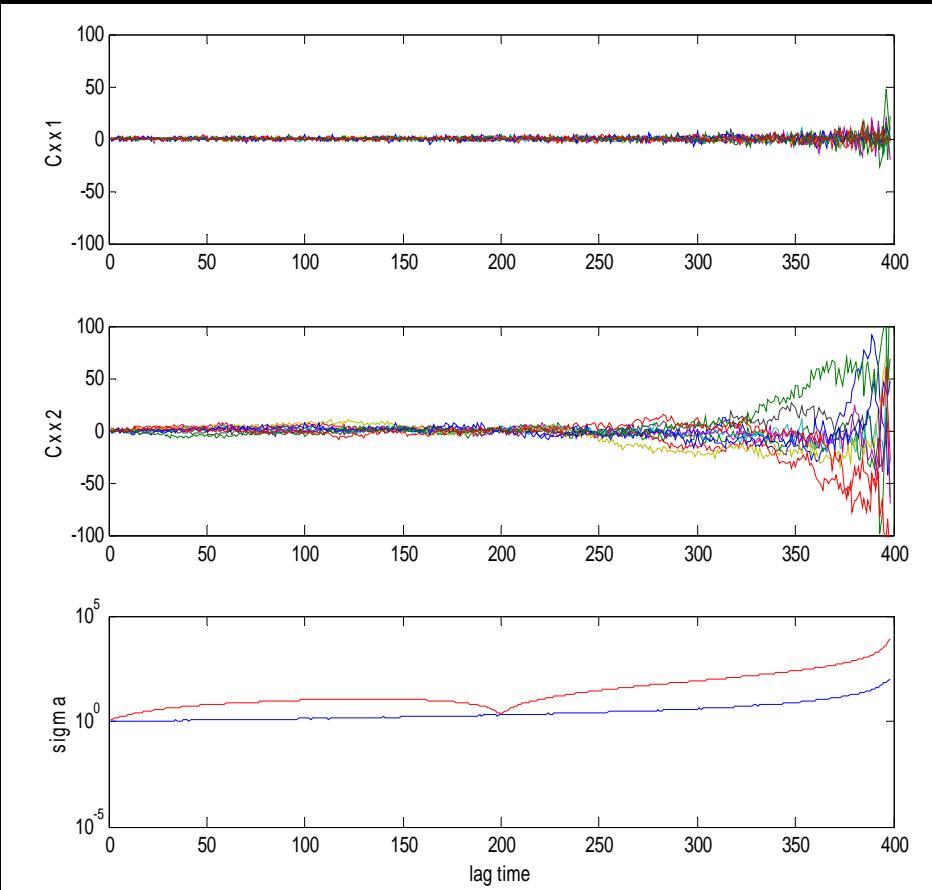
$$Z_N = \left| \frac{1}{n_t} \sum_{j_s=1}^{n_t} \bar{C}_{xx}^{\text{FT}}(j_s) \right| > c_{1-\alpha},$$

$$c_{1-\alpha} = \sqrt{2} \operatorname{Erfc}^{-1}(\alpha) \bar{\sigma}_{xx}^{\text{FT}}$$

E. coli adenylate kinase labeled w. Alexa fluor 555



The dip is counterintuitive

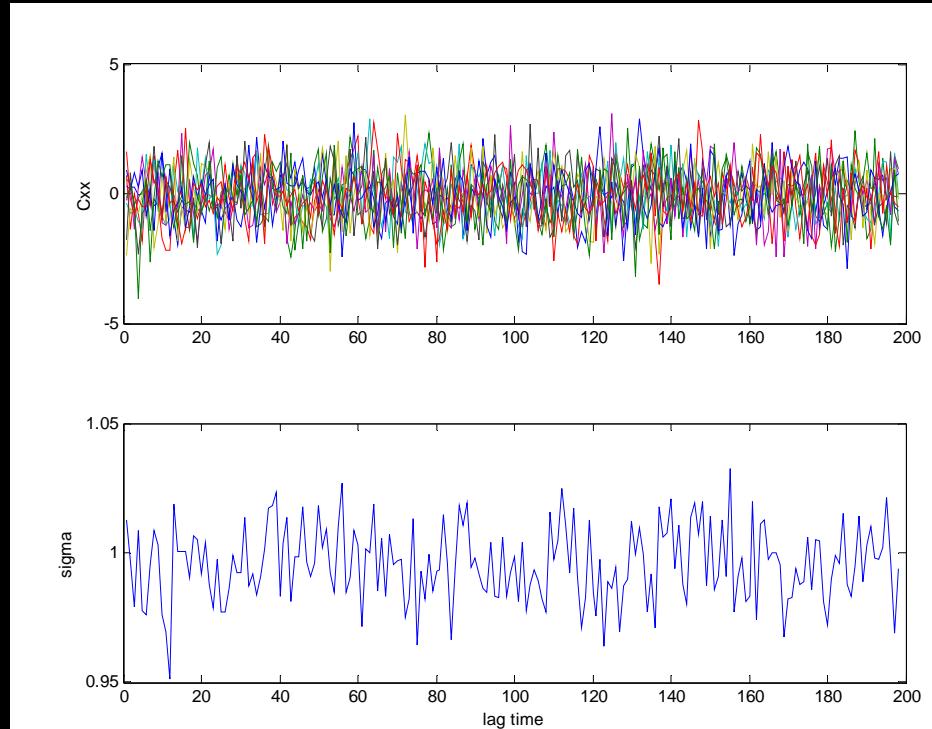


$$\bar{C}_{xx}(s) = \frac{1}{N - j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} - \bar{x}^2$$

$$\bar{C}_{xx}(s) = \frac{1}{N - j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2$$

$$\begin{aligned} \bar{C}_{xx}(s) = & \frac{1}{N - j_s} \sum_{i=1}^{N-j_s} x_i x_{i+j_s} \\ & - \left(\frac{1}{N - j_s} \sum_{i=1}^{N-j_s} x_i \right) \left(\frac{1}{N - j_s} \sum_{i=j_s+1}^N x_i \right) \end{aligned}$$

$C^{FT}(s)$



Fourier transform of the correation function

$$\bar{C}_{xx}^{FT}(s) = \frac{1}{N} \sum_{i=1}^N x_i x_{i+j_s} - \bar{x}^2$$

