

# Statistical analysis of time series: Gibbs measures and chains with complete connections

Rui Vilela Mendes

# Statistical analysis of time series data

- Time series

$$\cdots X_{-2} X_{-1} X_0 X_1 X_2 \cdots$$

$X \in Y$  : the *state space*

$Y^{\mathbb{Z}}$  : the *path space*

# Statistical analysis of time series data

- Time series

$$\cdots X_{-2} X_{-1} X_0 X_1 X_2 \cdots$$

$X \in Y$  : the *state space*

$Y^{\mathbb{Z}}$  : the *path space*

- **Statistical properties: (3 levels)**

- (1) Expectation values of the observables
- (2) Probability measures on state space  $Y$
- (3) Probability measures on path space  $Y^{\mathbb{Z}}$

# Statistical analysis of time series data

- Time series

$$\cdots X_{-2} X_{-1} X_0 X_1 X_2 \cdots$$

$X \in Y$  : the *state space*

$Y^{\mathbb{Z}}$  : the *path space*

- **Statistical properties: (3 levels)**

(1) Expectation values of the observables

(2) Probability measures on state space  $Y$

(3) Probability measures on path space  $Y^{\mathbb{Z}}$

- *Level 1, 2 and 3- statistical indicators.*

(*Mean partial sums, empirical measures (pdf's) and empirical process*)

# Statistical analysis of time series data

- Time series

$$\cdots X_{-2} X_{-1} X_0 X_1 X_2 \cdots$$

$X \in Y$  : the *state space*

$Y^{\mathbb{Z}}$  : the *path space*

- **Statistical properties: (3 levels)**

- (1) Expectation values of the observables
- (2) Probability measures on state space  $Y$
- (3) Probability measures on path space  $Y^{\mathbb{Z}}$

- *Level 1, 2 and 3- statistical indicators.*

(*Mean partial sums, empirical measures (pdf's) and empirical process*)

- **Analysis and reconstruction of the process:**

Purpose: To extract *Grammar* and *measure*

# Statistical analysis of time series data

- Time series

$$\cdots X_{-2} X_{-1} X_0 X_1 X_2 \cdots$$

$X \in Y$  : the *state space*

$Y^{\mathbb{Z}}$  : the *path space*

- **Statistical properties: (3 levels)**

(1) Expectation values of the observables

(2) Probability measures on state space  $Y$

(3) Probability measures on path space  $Y^{\mathbb{Z}}$

- *Level 1, 2 and 3- statistical indicators.*

(*Mean partial sums, empirical measures (pdf's) and empirical process*)

- **Analysis and reconstruction of the process:**

Purpose: To extract *Grammar* and *measure*

- Examples:

Hydrodynamic turbulence

Market fluctuations:

(there are analogies but the statistical indicators are different)

# Statistical analysis of time series data

- **Working hypothesis:** statistical methods are an appropriate tool to describe and reconstruct the market fluctuation process  
Related to modern view of the *efficient market* (expected value of abnormal returns is zero - Fama)  
Opposite view: behavioral component must always be included  
However: Behavioral trends not inconsistent with statistical description if the different reaction times of the market components as well as secondary reactions are taken into account (Olsen et al.)

# Statistical analysis of time series data

- **Working hypothesis:** statistical methods are an appropriate tool to describe and reconstruct the market fluctuation process  
Related to modern view of the *efficient market* (expected value of abnormal returns is zero - Fama)  
Opposite view: behavioral component must always be included  
However: Behavioral trends not inconsistent with statistical description if the different reaction times of the market components as well as secondary reactions are taken into account (Olsen et al.)
- Application of statistical tools requires:
  - (i) Stationary or asymptotically stationary process
  - (ii) Typical samples

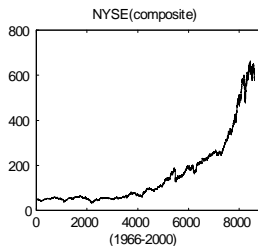
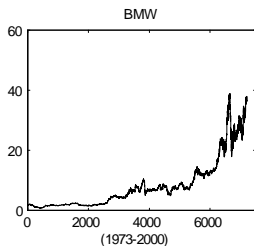
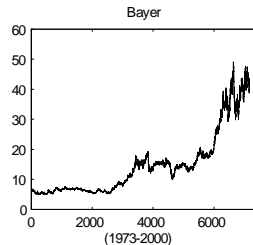
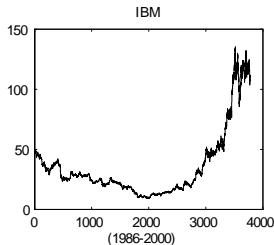


# Statistical analysis of time series data

- **Working hypothesis:** statistical methods are an appropriate tool to describe and reconstruct the market fluctuation process  
Related to modern view of the *efficient market* (expected value of abnormal returns is zero - Fama)  
Opposite view: behavioral component must always be included  
However: Behavioral trends not inconsistent with statistical description if the different reaction times of the market components as well as secondary reactions are taken into account (Olsen et al.)
- Application of statistical tools requires:
  - (i) Stationary or asymptotically stationary process
  - (ii) Typical samples
- Stocks as experimental probes revealing the mechanisms of the market process
  - (i)  $\Rightarrow$  preprocessing of the dataHigh-frequency versus low-frequency data  
Complexity versus statistics trade-off

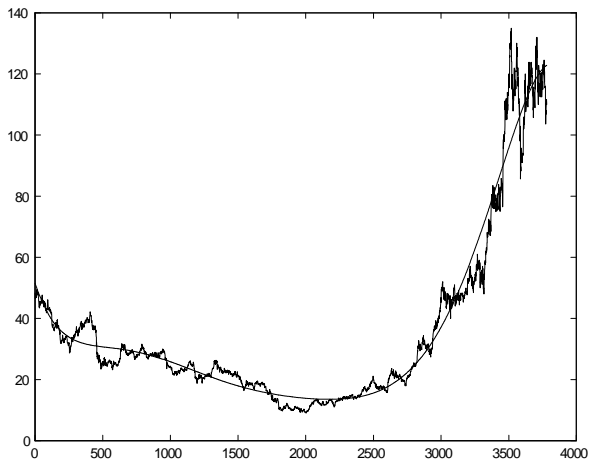
# Statistical analysis of time series data

Daily data  $p(t)$



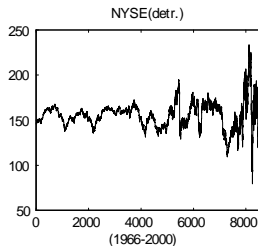
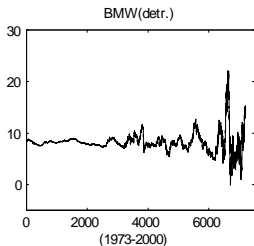
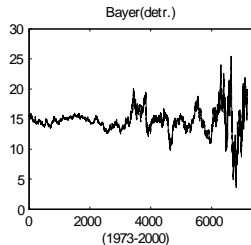
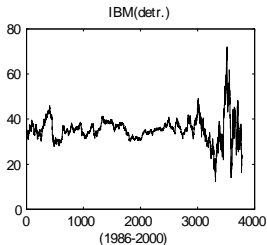
# Statistical analysis of time series data

## Detrending by a polynomial



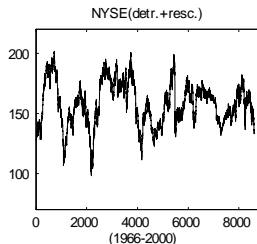
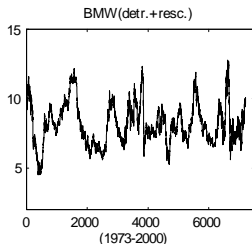
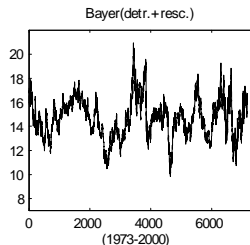
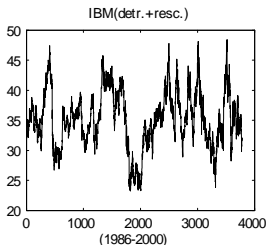
# Statistical analysis of time series data

Detrended data  $p(t) - q(t)$



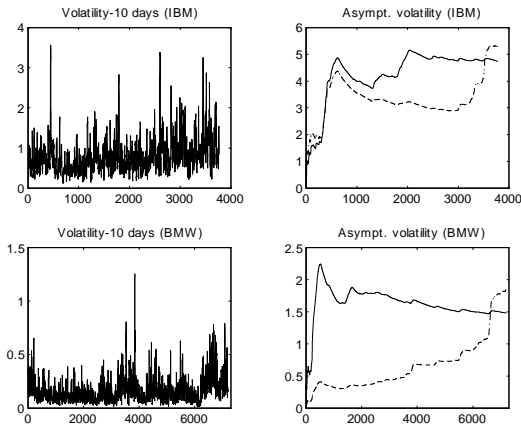
# Statistical analysis of time series data

Detrended and rescaled data  $x(t) = (p(t) - q(t)) \frac{\langle p(t) \rangle}{q(t)}$



# Statistical analysis of time series data

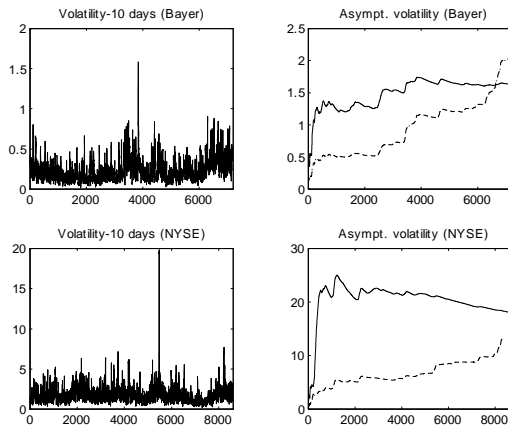
Ten-days window volatility and comparison of asymptotic volatility for the rescaled and non-rescaled data (IBM and BMW)



Local non-stationarity versus asymptotic stationarity

# Statistical analysis of time series data

Ten-days window volatility and comparison of asymptotic volatility for the rescaled and non-rescaled data (Bayer and NYSE)



Also direct test of stationarity computing the entropies of multi-symbol words in the first and second half of the samples

- $n$ —days return

$$r(t, n) = \log p(t + n) - \log p(t)$$



- $n$ —days return

$$r(t, n) = \log p(t + n) - \log p(t)$$

- (i) *Maximum* (over  $t$ ) of  $r(t, n)$

$$\delta(n) = \max_t \{r(t, n)\}$$

- $n$ —days return

$$r(t, n) = \log p(t + n) - \log p(t)$$

- (i) *Maximum* (over  $t$ ) of  $r(t, n)$

$$\delta(n) = \max_t \{r(t, n)\}$$

- (ii) *Moments* of the distribution of  $|r(t, n)|$

$$S_q(n) = \langle |r(t, n)|^q \rangle$$

# Statistical indicators

- $n$ —days return

$$r(t, n) = \log p(t + n) - \log p(t)$$

- (i) *Maximum* (over  $t$ ) of  $r(t, n)$

$$\delta(n) = \max_t \{r(t, n)\}$$

- (ii) *Moments* of the distribution of  $|r(t, n)|$

$$S_q(n) = \langle |r(t, n)|^q \rangle$$

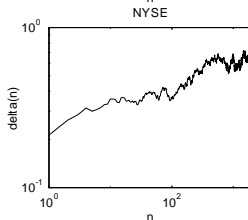
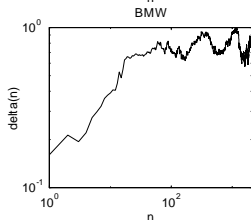
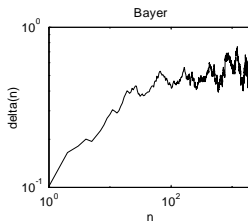
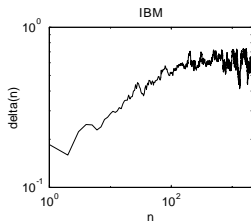
- (iii) If in some range ( $n = 2$  to  $n = 60$ )

$$S_q(n) \sim n^{\chi(q)}$$

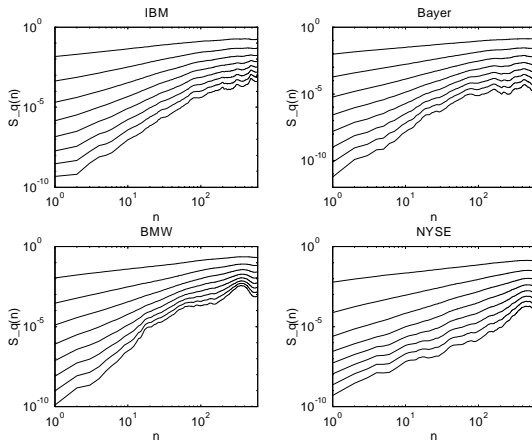
$\chi(q)$  is the *scaling exponent*

# Statistical indicators

Maximum  $\delta(n)$  of log-prices differences

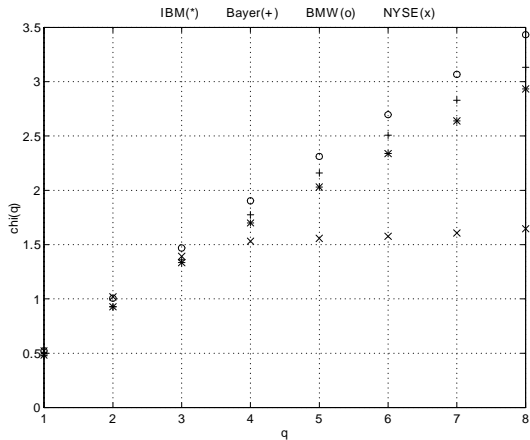


Moments of the  $|r(t, n)|$  distribution



# Statistical indicators

## Scaling exponent $\chi(q)$



- Main conclusions:
  - (a)  $\delta(n)$  is log-concave and probably asymptotically constant for large  $r$
  - (b)  $S_q(n)$  is a log-concave function of  $n$  with an inertial range
  - (c) The scaling law  $\chi(q)$  is an increasing concave function of  $q$
  - (d)  $\chi(1)$  in the scaling region ( $n = 2$  to  $n = 60$ ) is close to 0.5
  - (e) Scaling properties of NYSE somewhat different from the others

- Main conclusions:
  - (a)  $\delta(n)$  is log-concave and probably asymptotically constant for large  $r$
  - (b)  $S_q(n)$  is a log-concave function of  $n$  with an inertial range
  - (c) The scaling law  $\chi(q)$  is an increasing concave function of  $q$
  - (d)  $\chi(1)$  in the scaling region ( $n = 2$  to  $n = 60$ ) is close to 0.5
  - (e) Scaling properties of NYSE somewhat different from the others
- Properties (a) to (c) are shared by the turbulence data, but with different values for the statistical indicators (in turbulence data  $\chi(1) = \frac{1}{3}$ , here  $\chi(1) \approx 0.5 \Rightarrow$  essentially uncorrelated signal for  $n \geq 2$ )



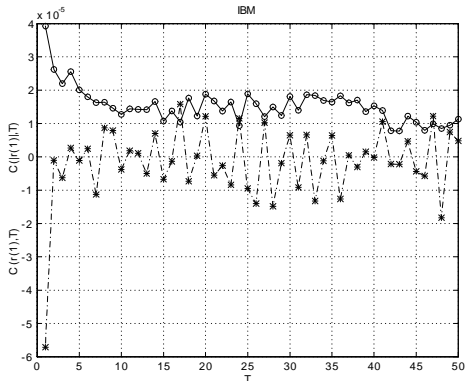
- Main conclusions:
  - (a)  $\delta(n)$  is log-concave and probably asymptotically constant for large  $r$
  - (b)  $S_q(n)$  is a log-concave function of  $n$  with an inertial range
  - (c) The scaling law  $\chi(q)$  is an increasing concave function of  $q$
  - (d)  $\chi(1)$  in the scaling region ( $n = 2$  to  $n = 60$ ) is close to 0.5
  - (e) Scaling properties of NYSE somewhat different from the others
- Properties (a) to (c) are shared by the turbulence data, but with different values for the statistical indicators (in turbulence data  $\chi(1) = \frac{1}{3}$ , here  $\chi(1) \approx 0.5 \Rightarrow$  essentially uncorrelated signal for  $n \geq 2$ )
- The behavior of the statistical indicators  $\delta(n)$ ,  $S_q(n)$  and  $\chi(q) \Rightarrow$  If the process is a topological Markov chain the transitions allowed by the transition matrix  $T$  must lie inside a strictly convex domain around the diagonal of  $T$

# Statistical indicators

Correlation function of one-day returns ( $\star$ ) and its absolute value ( $\circ$ )

$$C(r(1), T) = \langle r(t+T, 1) r(t, 1) \rangle$$

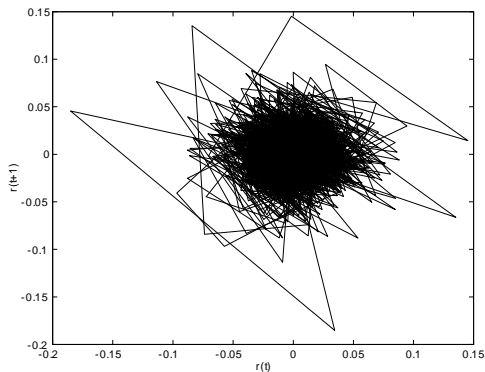
$$C(|r(1)|, T) = \langle |r(t+T, 1)| |r(t, 1)| \rangle$$



# Statistical indicators

## Dynamics of one-day returns

$$r(t, 1) \rightarrow r(t+1, 1)$$



## Why Gibbs measures are "natural"

- Events  $\{X_i\}$   
 $p_i =$  probability of event  $X_i$

## Why Gibbs measures are "natural"

- Events  $\{X_i\}$   
 $p_i$  = probability of event  $X_i$
- Constraints:
  - (i) Normalization,  $\sum_i p_i = 1$
  - (ii) Expectation value of known observables,  $\sum_i p_i F_k(X_i) = C_k$

# Looking for a Gibbs measure

## Why Gibbs measures are "natural"

- Events  $\{X_i\}$   
 $p_i$  = probability of event  $X_i$
- Constraints:
  - (i) Normalization,  $\sum_i p_i = 1$
  - (ii) Expectation value of known observables,  $\sum_i p_i F_k(X_i) = C_k$
- Maximum entropy principle (to maximize the uncertainty about what is not known) = use the most unbiased estimate

$$S = - \sum_i p_i \log p_i + \lambda_0 \sum_i p_i + \sum_k \lambda_k \sum_i p_i F_k(X_i)$$

$$\frac{\partial S}{\partial p_i} = 0 \quad \Rightarrow \quad -\log p_i - 1 + \lambda_0 + \sum_k \lambda_k F_k(X_i) = 0$$

$$p_i = \exp \left( -1 + \lambda_0 + \sum_k \lambda_k F_k(X_i) \right)$$

with  $\lambda_0, \lambda_1, \dots$  obtained from the constraints

# Looking for a Gibbs measure

- Coding by a finite alphabet  $\Sigma$

Space  $\Omega$  of orbits  $\omega = i_1 i_2 \cdots i_k \cdots, i_k \in \Sigma$

Dynamical law: a shift  $\sigma$

$$\sigma\omega = i_2 \cdots i_k \cdots$$

# Looking for a Gibbs measure

- Coding by a finite alphabet  $\Sigma$

Space  $\Omega$  of orbits  $\omega = i_1 i_2 \cdots i_k \cdots, i_k \in \Sigma$

Dynamical law: a shift  $\sigma$

$$\sigma\omega = i_2 \cdots i_k \cdots$$

- Grammar: set of allowed sequences in  $\Omega$



# Looking for a Gibbs measure

- Coding by a finite alphabet  $\Sigma$

Space  $\Omega$  of orbits  $\omega = i_1 i_2 \cdots i_k \cdots, i_k \in \Sigma$

Dynamical law: a shift  $\sigma$

$$\sigma\omega = i_2 \cdots i_k \cdots$$

- Grammar: set of allowed sequences in  $\Omega$
  - Sequences which coincide on the first  $n$  symbols:  $n$ -cylinder (or  $n$ -block) denoted by  $[i_1 i_2 \cdots i_n]$
- Probability measures over the cylinders

# Looking for a Gibbs measure

- Coding by a finite alphabet  $\Sigma$

Space  $\Omega$  of orbits  $\omega = i_1 i_2 \cdots i_k \cdots, i_k \in \Sigma$

Dynamical law: a shift  $\sigma$

$$\sigma\omega = i_2 \cdots i_k \cdots$$

- Grammar: set of allowed sequences in  $\Omega$
- Sequences which coincide on the first  $n$  symbols:  $n$ -cylinder (or  $n$ -block) denoted by  $[i_1 i_2 \cdots i_n]$   
Probability measures over the cylinders

- **Gibbs measure**

$$c_1 \leq \frac{\mu([i_1(\omega)i_2(\omega)\cdots i_n(\omega)])}{\exp(-nP + (S_n\phi)(\omega))} \leq c_2$$

$(S_n\phi)(\omega) = \sum_{k=0}^{n-1} \phi(\sigma^k\omega)$ ,  $\phi$  being Hölder continuous function on  $\Omega$  (the *potential*)

$P(\phi, G)$ : a function depending on potential and grammar (the *pressure of  $\phi$* )

# Looking for a Gibbs measure

- Relation to the entropy

$$h(\mu) = \lim_{n \rightarrow \infty} \frac{H_n}{n} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i_1 \dots i_n} \mu([i_1 i_2 \dots i_n]) \log \mu([i_1 i_2 \dots i_n])$$

# Looking for a Gibbs measure

- Relation to the entropy

$$h(\mu) = \lim_{n \rightarrow \infty} \frac{H_n}{n} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i_1 \dots i_n} \mu([i_1 i_2 \dots i_n]) \log \mu([i_1 i_2 \dots i_n])$$

- Variational principle: for each potential and grammar,  $\sup_{\eta} \{h(\eta) + \int \phi d\eta\}$  over all  $\sigma$ -invariant measures  $\eta$  is reached only for the Gibbs measure  $\mu$  and equals the pressure

$$P(\phi, G) = h(\mu) + \int \phi d\mu$$

# Looking for a Gibbs measure

- Relation to the entropy

$$h(\mu) = \lim_{n \rightarrow \infty} \frac{H_n}{n} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i_1 \dots i_n} \mu([i_1 i_2 \dots i_n]) \log \mu([i_1 i_2 \dots i_n])$$

- Variational principle: for each potential and grammar,  $\sup_{\eta} \{h(\eta) + \int \phi d\eta\}$  over all  $\sigma$ -invariant measures  $\eta$  is reached only for the Gibbs measure  $\mu$  and equals the pressure

$$P(\phi, G) = h(\mu) + \int \phi d\mu$$

- Potential may be chosen such that  $P = 0$  (*normalized potential*).  
Then

$$\phi(\omega) = \lim_{n \rightarrow \infty} \log \frac{\mu([i_1(\omega) \dots i_n(\omega)])}{\mu([i_2(\omega) \dots i_n(\omega)])}$$

(Practical use hindered by poor statistics of large blocks)

# Looking for a Gibbs measure

- **Gibbs measures for finite range potentials**

(finite range potentials approximate uniformly any Hölder continuous potential)

# Looking for a Gibbs measure

- **Gibbs measures for finite range potentials**

(finite range potentials approximate uniformly any Hölder continuous potential)

- Property of range- $r$  potentials: for all values  $i_1 i_2 \cdots i_n$  with  $n \geq r$

$$\mu([i_1 \cdots i_n]) = \frac{\mu([i_1 \cdots i_r]) \times \cdots \times \mu([i_{n-r+1} \cdots i_n])}{\mu([i_2 \cdots i_r]) \times \cdots \times \mu([i_{n-r+1} \cdots i_{n-1}])} \quad (1)$$

$\Rightarrow$

$$h(\mu) = - \sum_{i_1 \cdots i_k} \mu([i_1 \cdots i_k]) \log \frac{\mu([i_1 \cdots i_k])}{\mu([i_1 \cdots i_{k-1}])} = H_k - H_{k-1}$$

for all  $k \geq r$  if  $r > 1$ . If  $r = 1$

$$h(\mu) = H_1 H_k = - \sum_{i_1 \cdots i_k} \mu([i_1 \cdots i_k]) \log \mu([i_1 \cdots i_k])$$

# Looking for a Gibbs measure

- **Gibbs measures for finite range potentials**

(finite range potentials approximate uniformly any Hölder continuous potential)

- Property of range- $r$  potentials: for all values  $i_1 i_2 \cdots i_n$  with  $n \geq r$

$$\mu([i_1 \cdots i_n]) = \frac{\mu([i_1 \cdots i_r]) \times \cdots \times \mu([i_{n-r+1} \cdots i_n])}{\mu([i_2 \cdots i_r]) \times \cdots \times \mu([i_{n-r+1} \cdots i_{n-1}])} \quad (1)$$

$\Rightarrow$

$$h(\mu) = - \sum_{i_1 \cdots i_k} \mu([i_1 \cdots i_k]) \log \frac{\mu([i_1 \cdots i_k])}{\mu([i_1 \cdots i_{k-1}])} = H_k - H_{k-1}$$

for all  $k \geq r$  if  $r > 1$ . If  $r = 1$

$$h(\mu) = H_1 H_k = - \sum_{i_1 \cdots i_k} \mu([i_1 \cdots i_k]) \log \mu([i_1 \cdots i_k])$$

- $\Rightarrow$  **criterion to find the range of the potential:** range of the potential found when  $H_k - H_{k-1}$  tends to a constant.

Once the range is found, the potential may be constructed from the empirical weights  $\tilde{\mu}([i_1 \cdots i_k])$ .



# Looking for a Gibbs measure

- Another consequence of (1) is that for  $k > r$

$$\mu([i_1 \cdots i_{k+1}]) = \frac{\mu([i_1 \cdots i_k]) \mu([i_2 \cdots i_{k+1}])}{\mu([i_2 \cdots i_k])}$$

# Looking for a Gibbs measure

- Another consequence of (1) is that for  $k > r$

$$\mu([i_1 \cdots i_{k+1}]) = \frac{\mu([i_1 \cdots i_k]) \mu([i_2 \cdots i_{k+1}])}{\mu([i_2 \cdots i_k])}$$

- **Application to the market fluctuations:**

Five-symbols code  $\Sigma = \{-2, -1, 0, 1, 2\}$  for

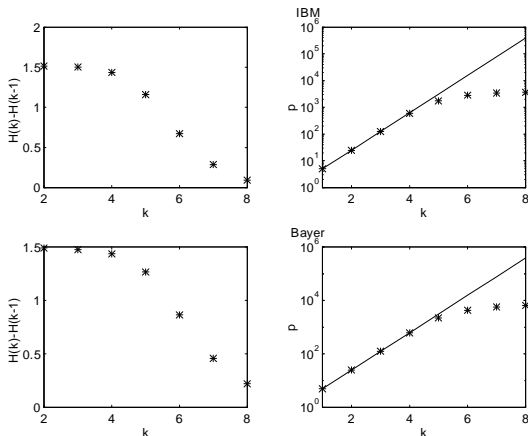
$$r(t) = \log p(t+1) - \log p(t)$$

Average  $\overline{r(t)}$  and standard deviation  $s = \sqrt{\overline{r^2(t)} - \overline{r(t)}^2}$

$$\begin{aligned} \left(r(t) - \overline{r(t)}\right) &> s && \iff 2 \\ s &\geq \left(r(t) - \overline{r(t)}\right) > \frac{s}{3} && \iff 1 \\ \frac{s}{3} &\geq \left(r(t) - \overline{r(t)}\right) > -\frac{s}{3} && \iff 0 \\ -\frac{s}{3} &\geq \left(r(t) - \overline{r(t)}\right) > -s && \iff -1 \\ -s &\geq \left(r(t) - \overline{r(t)}\right) && \iff -2 \end{aligned}$$

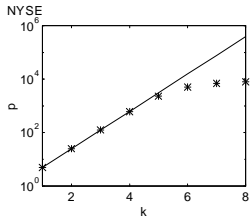
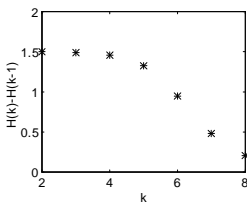
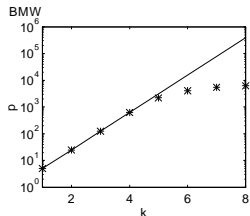
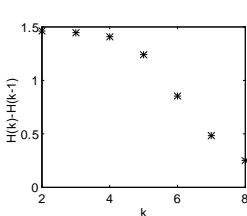
# Looking for a Gibbs measure

$H_k - H_{k-1}$  and the number of occurring blocks of size  $k$  (IBM and Bayer)



# Looking for a Gibbs measure

$H_k - H_{k-1}$  and the number of occurring blocks of size  $k$  (BMW and NYSE)



# Looking for a Gibbs measure

- Behavior of  $H_k - H_{k-1}$  quite different from hydrodynamic turbulence data

# Looking for a Gibbs measure

- Behavior of  $H_k - H_{k-1}$  quite different from hydrodynamic turbulence data
- To check whether a short-range potential is reliable :  
For successively higher  $k$  estimate

$$\mu_e([i_1 \cdots i_{k+1}]) = \frac{\tilde{\mu}([i_1 \cdots i_k]) \tilde{\mu}([i_2 \cdots i_{k+1}])}{\tilde{\mu}([i_2 \cdots i_k])}$$

then compare with the empirically observed  $\tilde{\mu}([i_1 \cdots i_{k+1}])$

# Looking for a Gibbs measure

- Behavior of  $H_k - H_{k-1}$  quite different from hydrodynamic turbulence data
- To check whether a short-range potential is reliable :  
For successively higher  $k$  estimate

$$\mu_e([i_1 \cdots i_{k+1}]) = \frac{\tilde{\mu}([i_1 \cdots i_k]) \tilde{\mu}([i_2 \cdots i_{k+1}])}{\tilde{\mu}([i_2 \cdots i_k])}$$

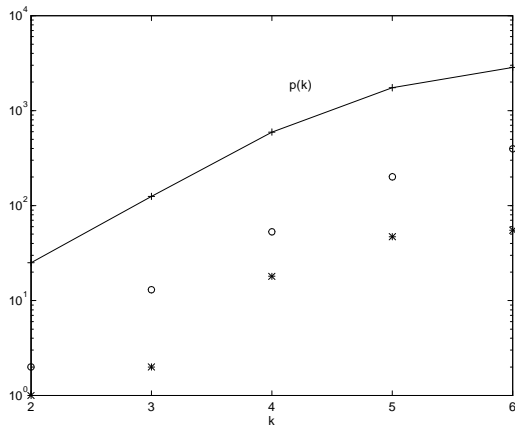
then compare with the empirically observed  $\tilde{\mu}([i_1 \cdots i_{k+1}])$

- Standard deviation of the relative positive errors

$$\varepsilon_k = \max \left( 0, \frac{\tilde{\mu}([i_1 \cdots i_{k+1}]) - \mu_e([i_1 \cdots i_{k+1}])}{\frac{1}{2} (\tilde{\mu}([i_1 \cdots i_{k+1}]) + \mu_e([i_1 \cdots i_{k+1}]))} \right)$$

# Looking for a Gibbs measure

Underestimation errors one (o) and two (\*) standard deviations away from the mean and the total number  $p(k)$  of observed blocks





# Looking for a Gibbs measure

- Large number of large deviation errors  
Correspond to blocks involving 2 and  $-2$   
Large deviations misrepresented by empirically constructed measure

# Looking for a Gibbs measure

- Large number of large deviation errors  
Correspond to blocks involving 2 and  $-2$   
Large deviations misrepresented by empirically constructed measure
- $?? \Rightarrow$  non-Gibbsian measure  
or  
 $?? \Rightarrow$  Gibbsian measure with long-range potential (sharp rise of  $H_k - H_{k-1}$  at  $k = 2$  followed by a very slow increase)

# Looking for a Gibbs measure

- Large number of large deviation errors  
Correspond to blocks involving 2 and  $-2$   
Large deviations misrepresented by empirically constructed measure
- $?? \Rightarrow$  non-Gibbsian measure  
or  
 $?? \Rightarrow$  Gibbsian measure with long-range potential (sharp rise of  $H_k - H_{k-1}$  at  $k = 2$  followed by a very slow increase)
- Large deviation analysis applied to the calculation of  $H_k$  consistent with this hypothesis

# Looking for a Gibbs measure

- Large number of large deviation errors  
Correspond to blocks involving 2 and  $-2$   
Large deviations misrepresented by empirically constructed measure
- $?? \Rightarrow$  non-Gibbsian measure  
or  
 $?? \Rightarrow$  Gibbsian measure with long-range potential (sharp rise of  $H_k - H_{k-1}$  at  $k = 2$  followed by a very slow increase)
- Large deviation analysis applied to the calculation of  $H_k$  consistent with this hypothesis
- In any case it needs an approach suited to deal with long-memory processes

# Chains with complete connections

- **Chain with complete connections (CCC)**

# Chains with complete connections

- **Chain with complete connections (CCC)**

①  $\forall a_i \in \Sigma$

$$P(X_1 = a_1, \dots, X_n = a_n) > 0$$

# Chains with complete connections

- **Chain with complete connections (CCC)**

- ①  $\forall a_i \in \Sigma$

$$P(X_1 = a_1, \dots, X_n = a_n) > 0$$

- ② The limit

$$\begin{aligned} & \lim_{m \rightarrow \infty} P(X_0 = a_0 | X_j = a_j, -m \leq j \leq -1) \\ &= P(X_0 = a_0 | X_j = a_j, j \leq -1) \end{aligned}$$

$$\text{exists } \forall a_i, j \leq -1$$

# Chains with complete connections

- **Chain with complete connections (CCC)**

- ①  $\forall a_i \in \Sigma$

$$P(X_1 = a_1, \dots, X_n = a_n) > 0$$

- ② The limit

$$\begin{aligned} & \lim_{m \rightarrow \infty} P(X_0 = a_0 | X_j = a_j, -m \leq j \leq -1) \\ &= P(X_0 = a_0 | X_j = a_j, j \leq -1) \end{aligned}$$

exists  $\forall a_i, j \leq -1$

- ③ There is a sequence  $(\gamma_m)_{m \geq 1}$  with  $\lim_{m \rightarrow \infty} \gamma_m = 0$ , such that for all  $\{a_j, b_j \in \Sigma, j \leq -1\}$  with  $a_j = b_j$  for  $-m \leq j \leq -1$

$$\left| \left( \frac{P(X_0 = a_0 | X_j = a_j, j \leq -1)}{P(X_0 = a_0 | X_j = b_j, j \leq -1)} - 1 \right) \right| \leq \gamma_m$$



# Looking for a Gibbs measure

- **Chain with complete connections and summable decay (CCCSd)**  
CCC with  $\sum \gamma_m < \infty$

# Looking for a Gibbs measure

- **Chain with complete connections and summable decay (CCCS)**  
CCC with  $\sum \gamma_m < \infty$
- Conditions 1. and 2. implicitly assumed for the pre-processed data  
Decays  $\gamma_m$  estimated from a typical sample of the process. From the empirical probabilities

$$P(a_0 | a_1 \cdots a_m A) = \frac{P(a_0 a_1 \cdots a_m A)}{P(a_1 \cdots a_m A)}$$

$A$  a block of arbitrary length

$$g(a_0 a_1 \cdots a_m) = \left( \frac{\max_A P(a_0 | a_1 \cdots a_m A)}{\min_A P(a_0 | a_1 \cdots a_m A)} - 1 \right)$$

$$\gamma_m = \max_{a_0 a_1 \cdots a_m} g(a_0 a_1 \cdots a_m)$$

If the statistics for long blocks is poor  $\Rightarrow$  large fluctuations in  $\gamma_m$

# Looking for a Gibbs measure

- **Chain with complete connections and summable decay (CCCSd)**  
CCC with  $\sum \gamma_m < \infty$
- Conditions 1. and 2. implicitly assumed for the pre-processed data  
Decays  $\gamma_m$  estimated from a typical sample of the process. From the empirical probabilities

$$P(a_0 | a_1 \cdots a_m A) = \frac{P(a_0 a_1 \cdots a_m A)}{P(a_1 \cdots a_m A)}$$

$A$  a block of arbitrary length

$$g(a_0 a_1 \cdots a_m) = \left( \frac{\max_A P(a_0 | a_1 \cdots a_m A)}{\min_A P(a_0 | a_1 \cdots a_m A)} - 1 \right)$$

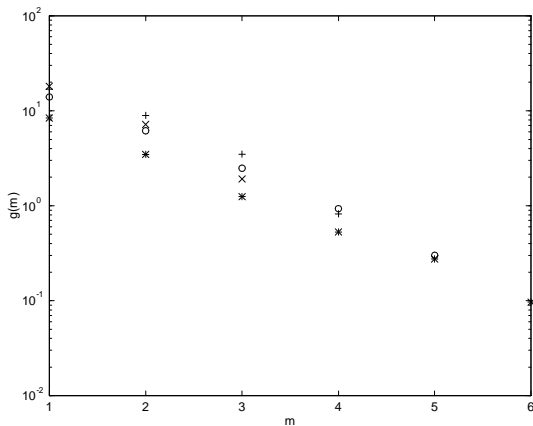
$$\gamma_m = \max_{a_0 a_1 \cdots a_m} g(a_0 a_1 \cdots a_m)$$

If the statistics for long blocks is poor  $\Rightarrow$  large fluctuations in  $\gamma_m$

- Better estimate of the decay behavior with  $\bar{g}(m) = \overline{g(a_0 a_1 \cdots a_m)}$ ,  
the average being taken over all sets  $a_0 a_1 \cdots a_m$  of size  $m$ .

# Chains with complete connections

$g(m)$  computed using  $A$  blocks of length 5 to 8 ( $\times, +, o, *$ )



# Chains with complete connections

- Result compatible with exponential decay  
⇒ summability of the  $\gamma_m$ 's

# Chains with complete connections

- Result compatible with exponential decay  
 $\Rightarrow$  summability of the  $\gamma_m$ 's
- **A CCC-process with summable decays is the  $\bar{d}$ -limit of its Markov approximations of order  $k$**

# Chains with complete connections

- Result compatible with exponential decay  
 $\Rightarrow$  summability of the  $\gamma_m$ 's
- **A CCC-process with summable decays is the  $\bar{d}$ -limit of its Markov approximations of order  $k$**
- $\bar{d}$ -distance

*Coupling* between two processes  $X = \{X_n\}$  and  $Y = \{Y_n\}$  is another process  $\{\tilde{X}_n, \tilde{Y}_n\}$  over  $\Sigma \times \Sigma$  such that the marginal probabilities of  $\tilde{X}$  and  $\tilde{Y}$  coincide with those of  $X$  and  $Y$

$$\bar{d}(X, Y) = \inf \left\{ P(\tilde{X}_0 \neq \tilde{Y}_0) : \left\{ \tilde{X}_n, \tilde{Y}_n \right\} \begin{array}{l} \text{is a stationary} \\ \text{coupling of } X \text{ and } Y \end{array} \right\}$$

# Chains with complete connections

- Result compatible with exponential decay  
 $\Rightarrow$  summability of the  $\gamma_m$ 's
- **A CCC-process with summable decays is the  $\bar{d}$ -limit of its Markov approximations of order  $k$**

- $\bar{d}$ -distance

*Coupling* between two processes  $X = \{X_n\}$  and  $Y = \{Y_n\}$  is another process  $\{\tilde{X}_n, \tilde{Y}_n\}$  over  $\Sigma \times \Sigma$  such that the marginal probabilities of  $\tilde{X}$  and  $\tilde{Y}$  coincide with those of  $X$  and  $Y$

$$\bar{d}(X, Y) = \inf \left\{ P(\tilde{X}_0 \neq \tilde{Y}_0) : \left\{ \tilde{X}_n, \tilde{Y}_n \right\} \text{ is a stationary coupling of } X \text{ and } Y \right\}$$

- $\bar{d}$ -distance tending to zero does mean that the processes will coincide after a certain time



# Chains with complete connections

- Result compatible with exponential decay  
 $\Rightarrow$  summability of the  $\gamma_m$ 's
- **A CCC-process with summable decays is the  $\bar{d}$ -limit of its Markov approximations of order  $k$**
- $\bar{d}$ -distance  
*Coupling* between two processes  $X = \{X_n\}$  and  $Y = \{Y_n\}$  is another process  $\{\tilde{X}_n, \tilde{Y}_n\}$  over  $\Sigma \times \Sigma$  such that the marginal probabilities of  $\tilde{X}$  and  $\tilde{Y}$  coincide with those of  $X$  and  $Y$   
$$\bar{d}(X, Y) = \inf \left\{ P(\tilde{X}_0 \neq \tilde{Y}_0) : \{\tilde{X}_n, \tilde{Y}_n\} \text{ is a stationary coupling of } X \text{ and } Y \right\}$$
- $\bar{d}$ -distance tending to zero does mean that the processes will coincide after a certain time
- **Perfect simulation** always understood in the  $\bar{d}$ -distance sense. It does not mean **perfect prediction** (Means that a process is constructed with the same conditional probabilities of the original one)

# Chains with complete connections

- **Simulation scheme by the sequence of canonical Markov approximations of finite order  $k$  ( $k$ -CMA)**

$k$ -CMA of a process  $X$  is a Markov chain  $Y^{(k)}$  of order  $k$  with conditional probabilities  $P^{(k)}$

$$P^{(k)}(a_0|a_1 \cdots a_k) = P(a_0|a_1 \cdots a_k) = \sum_A P(a_0|a_1 \cdots a_k A)$$

# Chains with complete connections

- **Simulation scheme by the sequence of canonical Markov approximations of finite order  $k$  ( $k$ -CMA)**

$k$ -CMA of a process  $X$  is a Markov chain  $Y^{(k)}$  of order  $k$  with conditional probabilities  $P^{(k)}$

$$P^{(k)}(a_0|a_1 \cdots a_k) = P(a_0|a_1 \cdots a_k) = \sum_A P(a_0|a_1 \cdots a_k A)$$

- For a CCC  $X$  with summable decays

$$\bar{d}(X, Y^{(k)}) \leq C\gamma_k \quad (2)$$

The property of the Markov approximation, essential for the approximation result (2), is

$$\inf_A P(a_0|a_1 \cdots a_k A) \leq P^{(k)}(a_0|a_1 \cdots a_k) \leq \sup_A P(a_0|a_1 \cdots a_k A) \quad (3)$$

meaning that for Markov approximation schemes, other than the canonical one, Eq.(2) holds provided (3) is satisfied

# Looking for a Gibbs measure

- **For the market fluctuation data:**

$\leq k$ —Markov approximation:

i) Empirical transition probabilities  $\tilde{P}(a_0|a_1 \cdots a_m)$  inferred from the probability of blocks of order  $m + 1$ . up to  $m_{Max}$

ii)  $\leq k$ —Simulation: look at the current block  $(a_1 \cdots a_k)$  and use the  $k$ —empirical probability to infer the next state  $a_0$ . If that block has not appeared in the training data, use the  $k - 1$  sized block  $a_2 \cdots a_k$  and the  $k - 1$  order empirical probabilities

# Looking for a Gibbs measure

- **For the market fluctuation data:**

- $\leq k$ —Markov approximation:

- i) Empirical transition probabilities  $\tilde{P}(a_0|a_1 \cdots a_m)$  inferred from the probability of blocks of order  $m+1$ . up to  $m_{Max}$

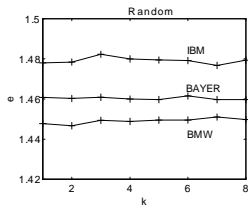
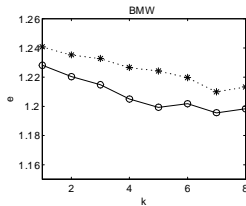
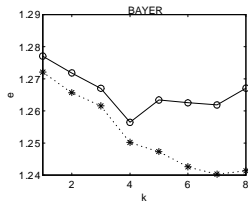
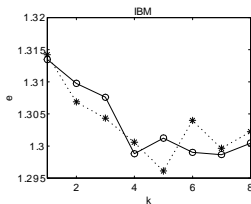
- ii)  $\leq k$ —Simulation: look at the current block  $(a_1 \cdots a_k)$  and use the  $k$ —empirical probability to infer the next state  $a_0$ . If that block has not appeared in the training data, use the  $k-1$  sized block  $a_2 \cdots a_k$  and the  $k-1$  order empirical probabilities

- Averaged squared error

$$e^2 = \left\langle (\tilde{a}_0 - a_0)^2 \right\rangle$$

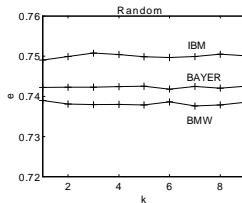
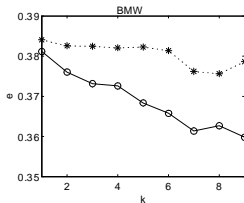
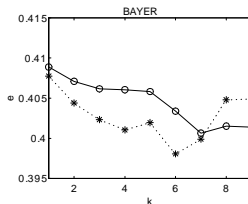
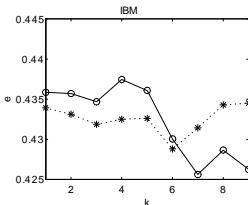
# Chains with complete connections

The past predicting the future (o) and the future predicting the past (\*), compared to random choice



# Chains with complete connections

The past predicting the future (o) and the future predicting the past (\*), compared to random choice



# Chains with complete connections

## Main conclusions:

- Average prediction better than random choice.
- Main improvement results from correct accounting of the two-symbol probabilities ( $k = 1$ )
- Small (but consistent) improvement using past information up  $k = 4$  or 5. No significant improvement using higher order approximations
- Bulk of data represented by a short-memory process. Nevertheless there is evidence for a small long-memory component that is captured by the higher-order Markov approximations
- There is a maximum  $k = k_m$  that should be used for the simulation process
- Inter-companies prediction: improvement coming from the one-symbol probabilities (as compared to random choice) is obtained. For the long-memory component the behavior is company-dependent.



# Statistical analysis of market data

## General conclusions

- 1 Bulk of the market fluctuation process is a short-memory process. In addition it has a small long-memory component associated with the large fluctuations of the returns.
- 2 Existence of the long-memory component suggests the *chains with complete connections and summable decays* as a framework
- 3 Although the decays may be exponentially converging, the lack of accurate data for long blocks prevents an accurate description by a finite range Gibbs potential.
- 4 The sequence of empirical based  $\leq k$ -Markov approximations seems the most unbiased simulation of the process. Eventual convergence in the  $\bar{d}$ -distance sense, because of summable decays.
- 5 For the future: high frequency market data. Beware of the possibly multi-scale and multi-component nature of the processes.

# References

- # J.-R. Chazottes, E. Floriani and R. Lima; *Relative entropy and identification of Gibbs measures in dynamical systems*, J. of Stat. Phys. 90 (1998) 697-725.
- # E. F. Fama; *Market efficiency, long-term returns and behavioral finance*, Journal of Financial Economics 49 (1998) 283-306.
- # R. B. Olsen, M. M. Dacorogna, U. A. Müller and O. V. Pictet; *Going back to basics - Rethinking market efficiency*, Olsen & Associates discussion paper RBO.1992-09-07.
- # D. Ruelle; *Thermodynamic formalism*, Encyclopedia of Mathematics and its Applications 5, Addison-Wesley 1978.
- # O. Onicescu and G. Mihoc; *Le comportement asymptotique des chaines à liaisons complètes*; Disq. Math. Phys. 1 (1940) 61-62.
- # M. Iosifescu and S. Grigorescu; *Dependence with complete connections and its applications*, Cambridge U. P., Cambridge 1990.
- # D. S. Ornstein and B. Weiss; *How sampling reveals a process*, The Annals of Prob. 18 (1990) 905-930.

- # X. Bressaud, R. Fernandez and A. Galves; *Speed of  $\bar{d}$ –convergence for Markov approximations of chains with complete connections. A coupling approach*, Stoch. Proc. and Appl. 83 (1999) 127-138.
- # F. Comets, R. Fernandez and P. A. Ferrari; *Processes with long memory: Regenerative construction and perfect simulation*, arXiv:math.PR/0009204.
- # RVM, R. Lima and T. Araújo; *A process-reconstruction analysis of market fluctuations*, Int. Journal of Theoretical and Applied Finance, 5 (2002) 797