

Master Thesis

Subject: <<Motifs of Network Models>>

Student : George Argyris

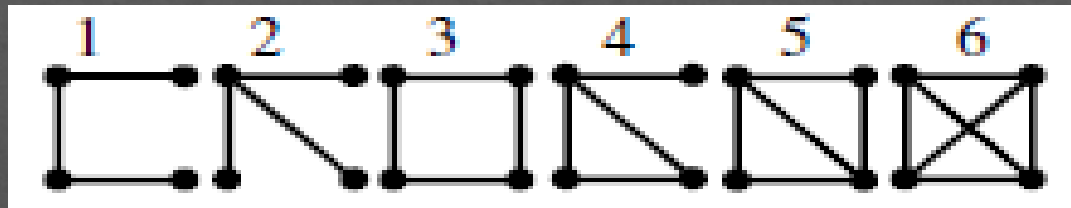
Advisor: Dimitris Kugiumtzis

AUTH Thessaloniki December 2016

THE GOAL: Separate the concepts of randomness, smallworldness, scalefreeness

Motifs

- ◇ Generally: statistically significant sub-graphs of a graph that occurs recurrently within the graph
- ◇ In our study, motifs are connected (tetrads) 4's nodes:



- ◇ Counting a graph's motifs:

-Work on the adjacency matrix

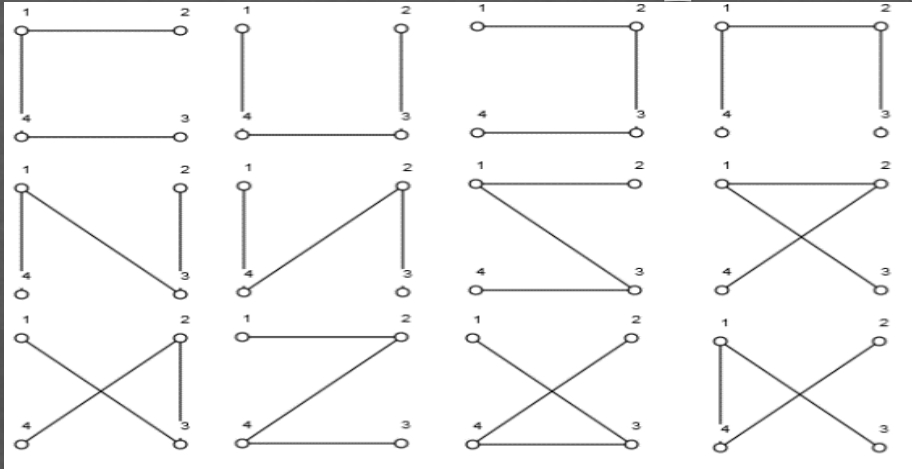
Examine all 4x4 sub-matrices of A ($\binom{n}{4} = \frac{n!}{4!(n-4)!} = \frac{n \cdot (n-1) \cdot (n-2) \cdot (n-3)}{24}$)

-Symmetric graphs' problem: The algorithm was computationally complex because one motif is depicted by more than one adjacency matrices.

Solution: Same degree distribution for graphs of size: $n \leq 4 \Leftrightarrow$ isomorphic motifs

-I used Brain Connectivity Toolbox in order to further reduce computational complexity.

Motif 1 frequency in ER random graphs

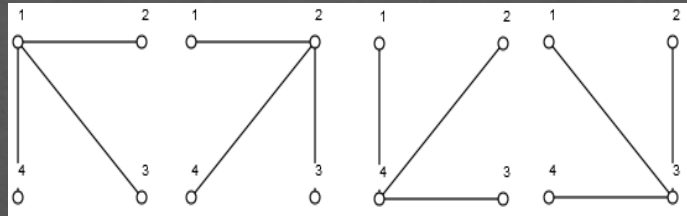


◇ 12 symmetries

◇ occurrence probability for each symmetry: $p^3 \cdot (1-p)^3$

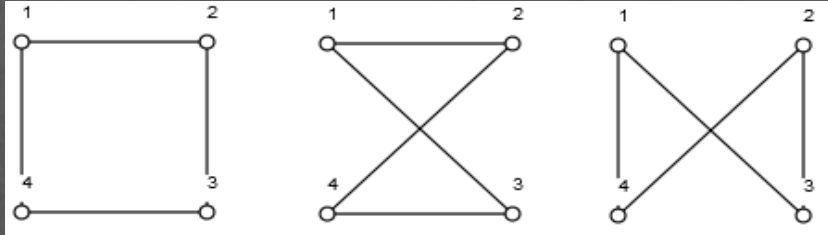
◇ Motif's 1 frequency: $f_1 = 12 \cdot \binom{n}{4} \cdot p^3 \cdot (1-p)^3$

Motif 2 frequency in E-R random graphs



- ◇ 4 symmetries
- ◇ Probability of occurrence of each one of them: $p^3 \cdot (1-p)^3$
- ◇ Motif's 2 frequency: $f_2 = 4 \cdot \binom{n}{4} \cdot p^3 \cdot (1-p)^3$

Motif 3 frequency in E-R random graphs



◇ 3 symmetries

◇ Probability of occurrence of each one of them: $p^4 \cdot (1-p)^2$

◇ Motif's 3 frequency: $f_{\downarrow 2} = 3 \cdot \binom{n}{4} \cdot p^4 \cdot (1-p)^2$

Motif frequency of Erdos-Renyi random graphs

- ◇ Frequency of motif i :

$$f_i(n,p) = N \binom{n}{m_i} p^l (1-p)^{6-l} \quad (1)$$

- ◇ Relative frequency of motif i :

$$F_i(n,p) = f_i(n,p) / \sum_{i=1}^6 f_i(n,p)$$

- ◇ Motifs' relative frequency vector:

$$F = (F_1, F_2, F_3, F_4, F_5, F_6)$$

Highlight 1

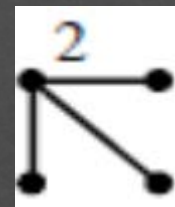
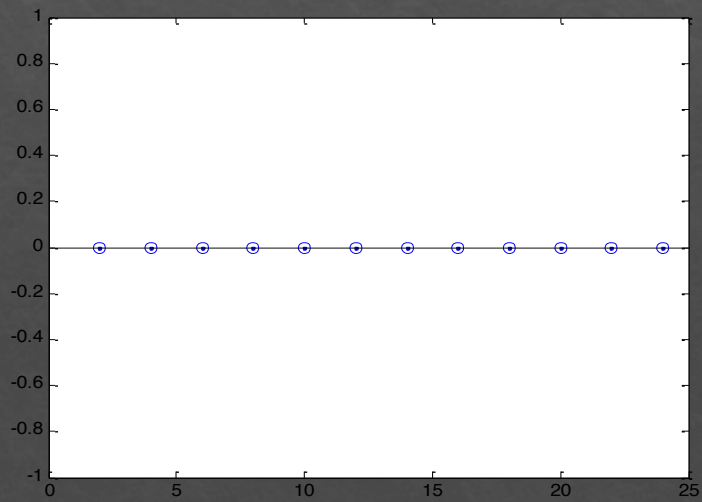
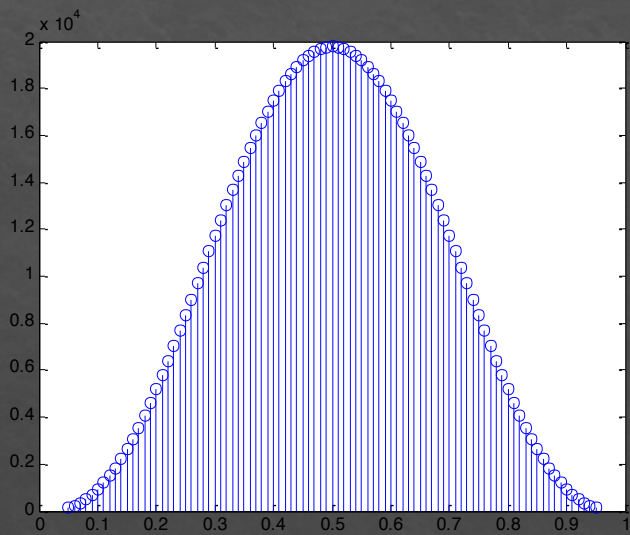
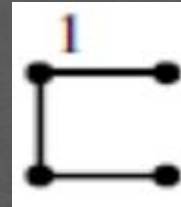
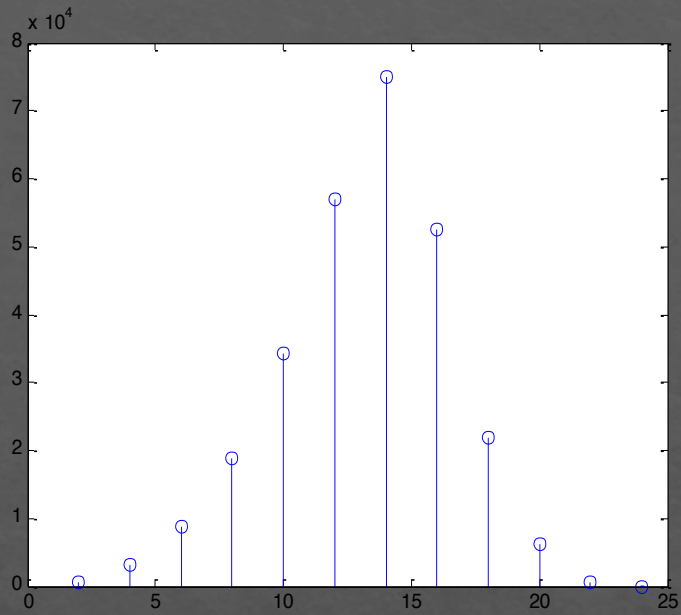
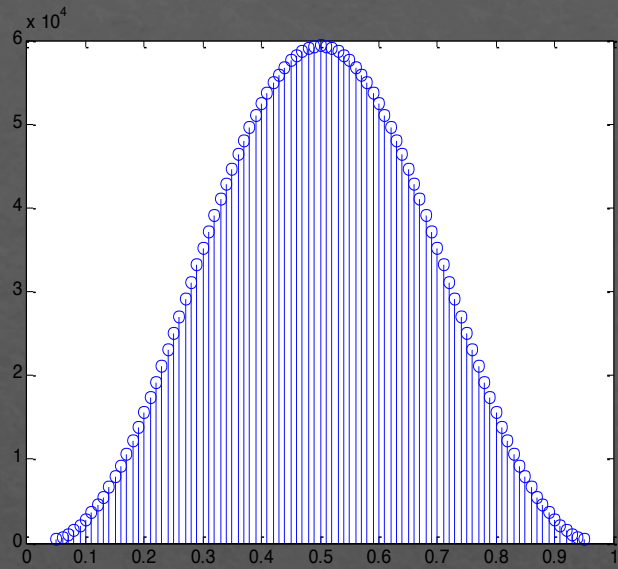
Motifs' frequency comparison between ER and WS networks

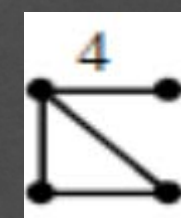
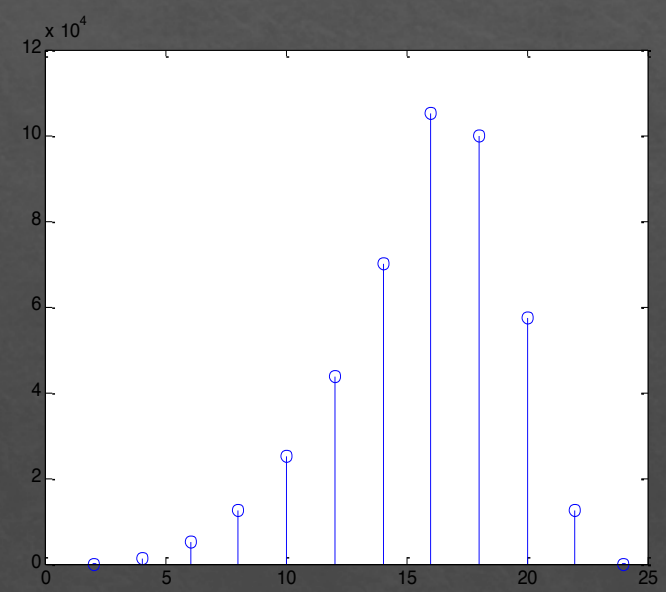
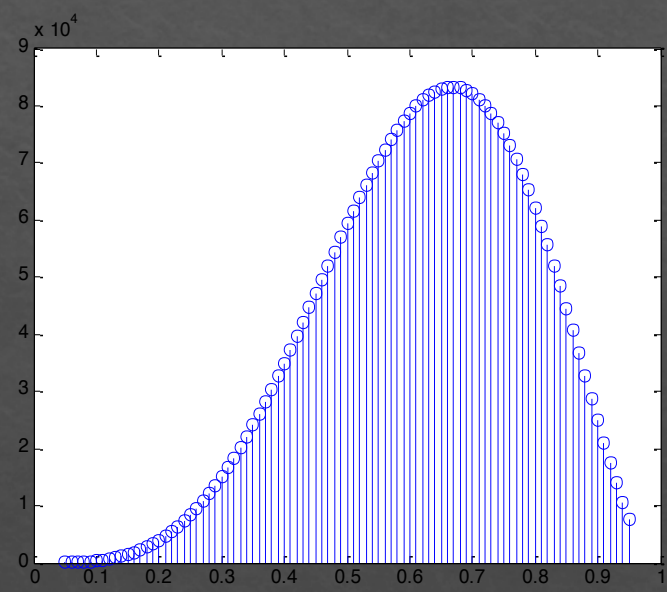
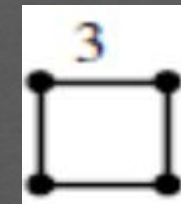
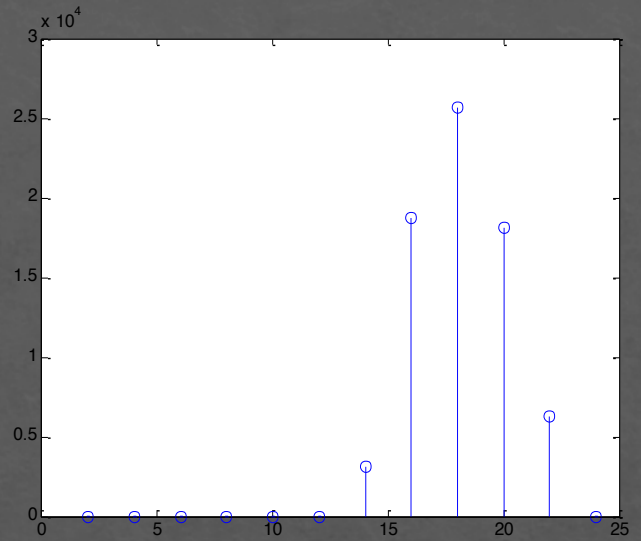
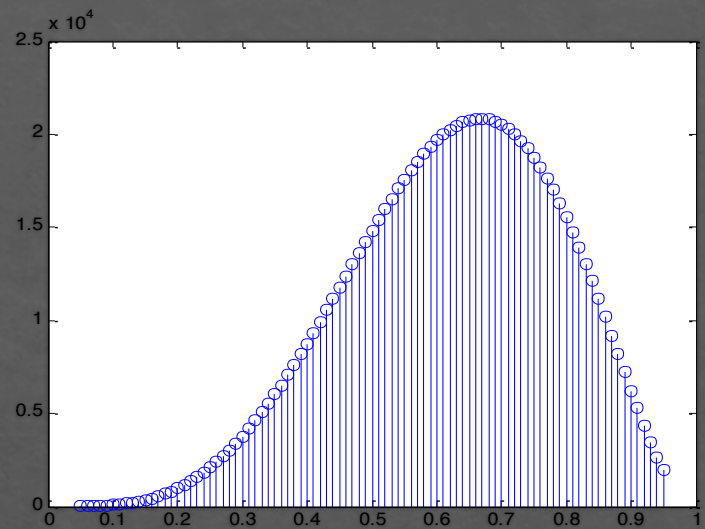
We counted the motifs' frequency for 25 networks of 25 nodes (random and small world ($p_d = 0$)).

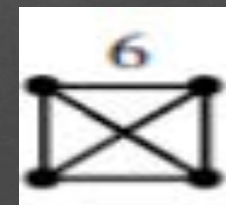
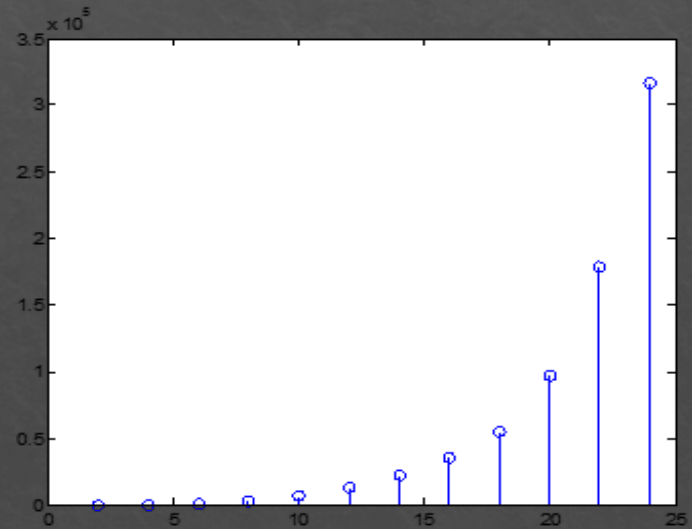
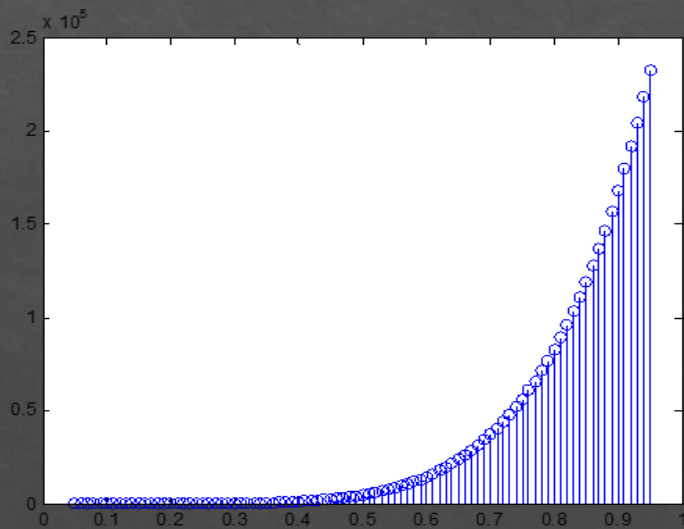
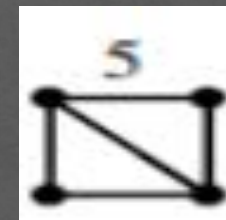
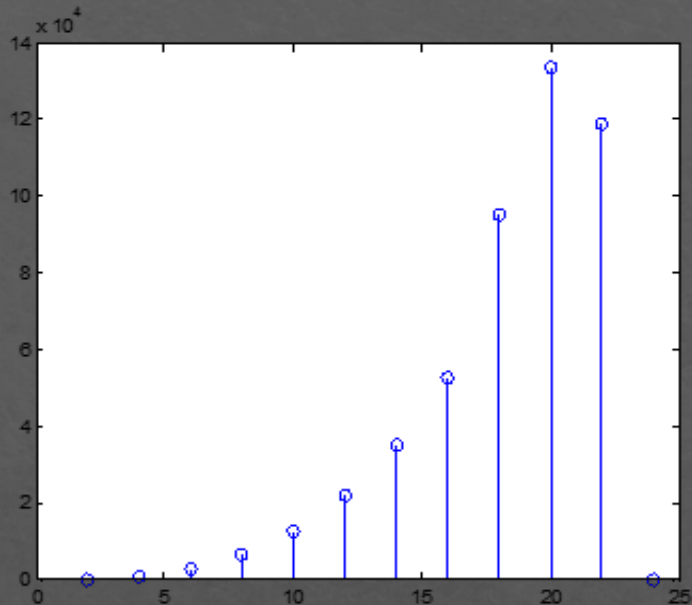
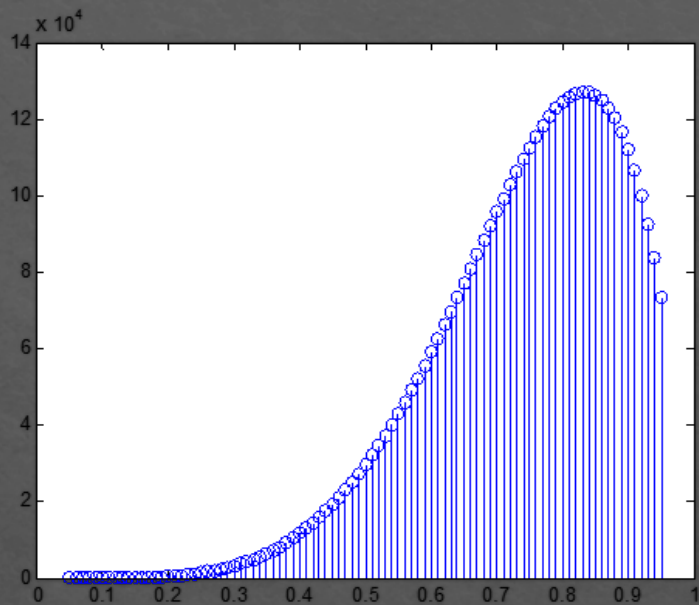
- ◇ In case of ER networks we constructed the graphs analytically, for different values of p , according to this function:

$$m_l(p) = 25 \cdot f_l(25, p) = 25 \cdot N_{mi} \cdot \binom{n-1}{l} \cdot p^l \cdot (1-p)^{n-l}$$

- ◇ In case of WS networks we constructed the graphs empirically, for different values of k and $p_d = 0$.







Highlights

- ◇ The smaller the euclidean distance between the motifs' relative frequency vector of an observed network and the expected motifs' relative frequency vector calculated previously the more random the network is.
- ◇ Maybe, we have same degree distributions for 5 of 6 motifs with different parameters independently the construction model.
- ◇ Motif 2 of great importance! As p of WS model approximates 1 the frequency of that motif tends to the expected value of a random one.

Problem-Method

- ◇ Suppose we have an observed network, could we find which model generates it? Or which model best approximates it?

Method

- ◇ Construction of the motifs' relative frequency vector for each one of the models and for each of their parameters as $n=25, 50, 75$.

E-R: $p=0.1, 0.2, \dots, 0.9$

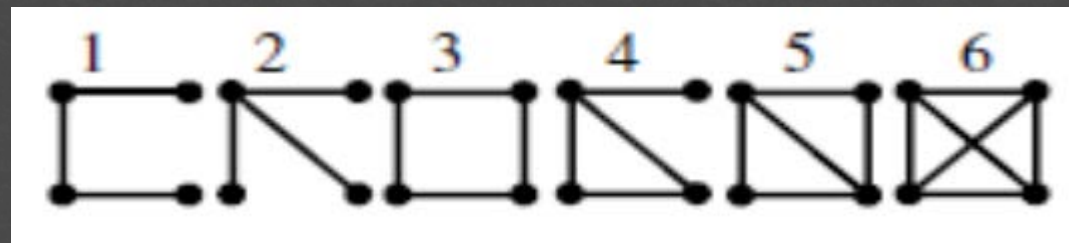
W-S: $k=2, 4, 6, \dots, \quad p_d = 0.1, 0.2, \dots, 0.9$

B-A: $m_i =$ each one of the initial motifs, $i=1, \dots, 6, \quad m_0=1, 2, 3, 4$

EXP: $\text{exp}=2.1, 2.2, 2.3, \dots, 2.9$

- ◇ Finding the corresponding vector of the observed network.
- ◇ The model that generates the network is the one that minimizes the euclidean distance.

- ◇ We expect intuitively that:



Why method does not work well

◇ Method does not work well because:

1. Networks constructed by one model are classified to close models.

Ex. ER(p) networks are mapped ER(p'), $p' = p \pm 0.1$

WS(k, p) are mapped to WS(k', p') with $p' = p \pm 0.1$ and $k' = k \pm 2$

2. Density is the most critical parameter that affects the relative frequency of a vector.

Solution

If parameter density of the observed network is known, we could calculate some of the parameters of the models. So, the models that generate an observed graph become less, and the method more accurate.

The parameter calculation

- ◆ Concerning density d from an observed network, we could calculate some of the parameters of the models as below:

Erdos-Renyi: $p = d$

Watts-Strogatz: $d = \text{GraphOrder} / \text{MaximumOrder} = n \cdot k / 2 / n \cdot (n-1) / 2 = k / (n-1) \Rightarrow k = d \cdot (n-1)$

Barabasi-Albert:

$$m \downarrow 0 = (d \cdot (n-1) \cdot n / 2 - 4) / n - 4$$

Success rate of the method for standard density

n=50

	BA.m1.1	BA.m1.2	BA.m1.3	BA.m1.4	BA.m2.1	BA.m2.2	BA.m2.3	BA.m2.4
RN	0	4	2	0	1	1	0	0
WS.0.0	0	0	0	0	0	0	0	0
WS.0.1	0	0	0	0	0	0	0	0
WS.0.2	0	0	0	0	0	0	0	0
WS.0.3	0	0	0	0	0	0	0	0
WS.0.4	0	0	0	0	0	0	0	0
WS.0.5	0	0	0	0	0	0	0	0
WS.0.6	0	0	0	0	0	0	0	0
WS.0.7	0	0	0	0	0	0	0	0
WS.0.8	0	0	0	0	0	0	0	0
WS.0.9	0	0	0	0	0	0	0	0
BA.m1	58	8	49	42	46	5	48	46
BA.m2	17	3	0	13	17	0	0	6
BA.m3	0	37	0	5	0	49	0	2
BA.m4	0	0	9	0	0	0	4	1
BA.m5	14	48	32	11	30	45	34	11
BA.m6	11	0	8	29	6	0	14	34

	BA.m3.1	BA.m3.2	BA.m3.3	BA.m3.4	BA.m4.1	BA.m4.2	BA.m4.3	BA.m4.4
RN	0	6	3	0	0	1	0	2
WS.0.0	0	0	0	0	0	0	0	0
WS.0.1	0	0	0	0	0	0	0	0
WS.0.2	0	0	0	0	0	0	0	0
WS.0.3	0	0	0	0	0	0	0	0
WS.0.4	0	0	0	0	0	0	0	0
WS.0.5	0	0	0	0	0	0	0	0
WS.0.6	0	0	0	0	0	0	0	0
WS.0.7	0	0	0	0	0	0	0	0
WS.0.8	0	0	0	0	0	0	0	0
WS.0.9	0	0	0	0	0	0	0	0
BA.m1	41	6	49	49	36	11	49	35
BA.m2	28	0	0	13	11	3	0	11
BA.m3	0	38	0	4	0	34	0	1
BA.m4	0	0	3	3	0	0	8	5
BA.m5	17	50	39	4	28	51	27	11
BA.m6	14	0	6	27	25	0	16	35

Application to real data

- ◇ The Morgan Stanley Capital International's (MSCI) dataset contains the capitalization indices for 55 markets.
- ◇ It contains 1305 daily returns for each market from 5 March 2004 to 5 March 2009.
- ◇ We separate the dataset into 87 and 29 subdatasets each one of them contains 15 and 45 returns respectively in chronological order. For each one of the subdatasets we have an association network which is constructed by the correlation of the time series.

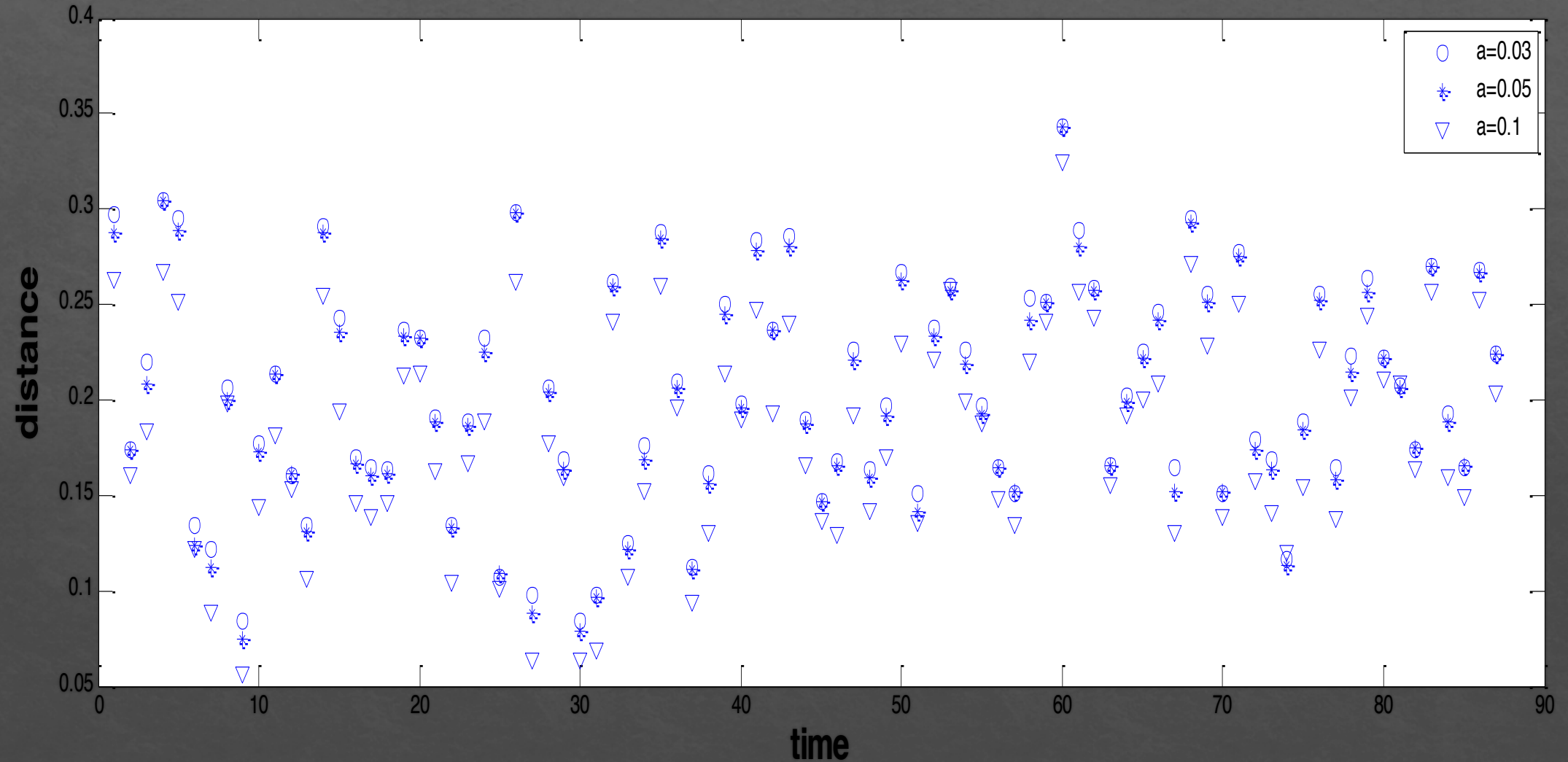
Statistical procedures used

For the statistical analysis of the multivariate time series we used:

- ◆ For each one of the subdataset and for each one of the time series a prewhiten method.
- ◆ A significance test for finding the statistically significant cross correlation. In that test there is a significance level α . As α gets large, the cross-correlation that are statistically significant are getting large. Also, α determines the density d of the association network.

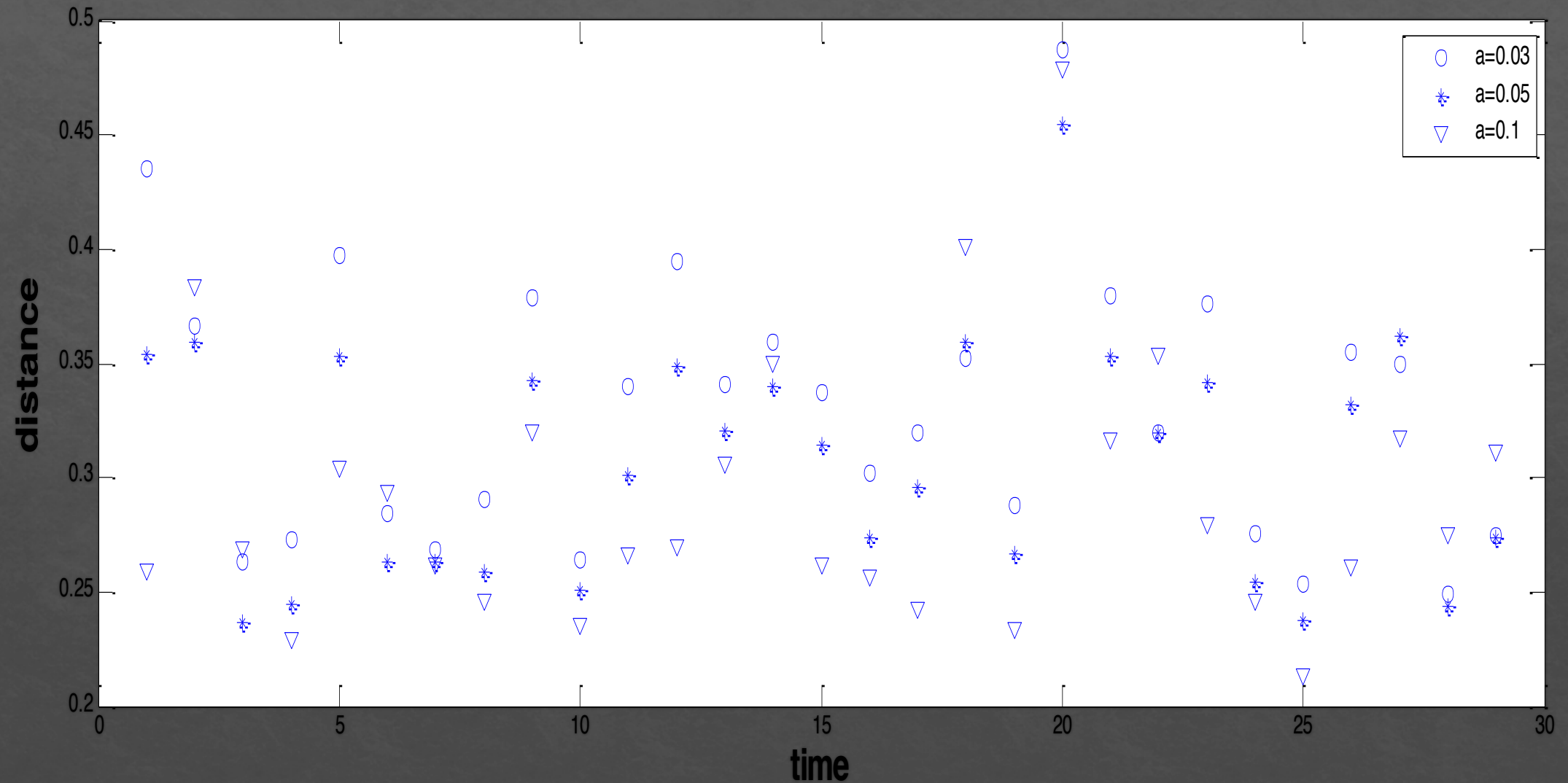
The randomness indicator for the association networks for the case of 87 networks

Distances from random networks of time series of networks



The randomness indicator for the association networks for the subdataset of 29 networks

Distances from random networks of time series of networks



Conclusion

Statement 1.

Let G be an $ER(p)$ graph and G_1, G_2 be subgraphs of same order and size.

The frequency of G_1 is higher than the frequency of G_2 if and only if G_1 has more symmetries than G_2 .

Statement 2.

Subgraph G achieves highest frequency in an $ER(p)$ graph when $p = d_G$.

Future research activities

- ◇ Find the motifs of a $WS(n, k, 0)$ model analytically.
- ◇ Find which motif triad better recognizes the model and make a 3D representation.
The statement 2 may help to this direction.
- ◇ Application to other multivariate time series.
- ◇ Instead of euclidean distance, we could train a perceptron which works if-f the data are linearly separable and we compare only 2 models.
- ◇ Same analysis for motifs of size 5 is feasible.
- ◇ Extendable to directed networks.
- ◇ The method might be extendable to other network models.