

Complexity and Fitness of Methylated DNA sites and RNAs in the context of Breast Cancer Tissues

June 22, 2023

1 Aim of the study

Highlighting the nestedness structure of the bipartite network representing correlations between methylated DNA sites and RNAs.

2 Methods and Data

2.1 Nestedness

Nested structures appear in the nature such that in ecological and socio-economic systems. In ecological systems, for instance, there is a nested distribution of species habitats. Where as some species only live in particular locations, some others live in a large variety of environments including those particular locations. In the context of network theory, a perfectly nested structure is such that for any pair of node (i, j) with $k_i < k_j$ their respective degree (number of neighbors), the set of nodes i is interacting with is included in the set of nodes j is interacting with. The adjacency matrix of such a network shows a triangular shape when columns and rows are sorted accordingly to the degree ordering of respective nodes. There are many metrics (other than degrees) permitting to infer the nested structure of a given network. Nestedness can appear in unipartite and bipartite networks. In the context of economical systems, especially the world trade network seen as a country-product bipartite network, recent research shows the existence of nestedness, see a complete review in [mariani19]. In the context of biological systems nestedness structures have also been studied. In [cantor17] different biological scales are considered. However, it seems that DNA-RNA interactions have not been studied yet. We choose the metric called fitness-complexity introduced in [tacchella12] to investigate the nestedness structure of DNA-RNA bipartite network. Fitness and complexity is a non-linear and iterative algorithm permitting to infer the nestedness in bipartite network such as economical networks. In this context, the fitness is a country quality, and complexity is a product quality. In our approach we try to make an analogy where RNAs play the role of products and DNAs the role of countries. Here are the economical definitions of fitness and complexity.

Fitness and Complexity metrics

$$\begin{cases} F_d^{(n)} = \sum_r B_{dr} Q_r^{(n-1)} \\ Q_r^{(n)} = \left(\frac{1}{\sum_d B_{dr} \frac{1}{(F_d^{(n-1)})^\gamma}} \right)^{1/\gamma} \end{cases} \quad (1)$$

with B_{dr} an element of the bi-adjacency matrix representing the correlation (positive or negative) interaction between a Methylated DNA d and an RNA r , $F^{(n)}$ and $Q^{(n)}$ the Fitness and Complexity vector measured at iteration n . We set γ to 2.

Product's complexity

2.2 Data

We used an open access cancer dataset from GDC Data Portal. From raw data consisting in methylated DNA and RNA's beta-values for up to 841 tumorous and normal breast tissues, we built a Matrix of Pearson Correlation Coefficients between Methylated DNA and RNA from these beta-values. A network representing the Pearson correlation coefficients between pairs (M-DNA,RNA) consists in a bipartite graph with coefficients as weight of the links containing. We also project DNAs into spacial clusters of length 40K base-pairs, this turns $N_{DNA} = 364285$ DNAs into a set of $N_{cl} = 93690$ clusters. The nature of correlations are either positive or negative.

We propose to measure complexity of RNAs from the complexity-fitness metric proposed in [tacchella12].

3 The case of ESR1 synthesis and regulation

The set of ESR1 synthesis and regulation consists in 508 RNAs and 4959 clusters of DNAs. We first construct the Pearson correlation matrix between clustered DNAs and RNA then we measure the complexity and fitness of nodes from the sub network related to the ESR1 regulation sub set.

Table 1: Highest Clustered DNAs and RNAs in terms of fitness and complexity scores, 30 RNAs and 228 Clusters, with $\rho_c = 0.7$. ¹: cg02730804, cg04077850, cg14800014, cg13049992, cg09168797 and cg01977473, ²: cg27121959, ³:cg03633268.

RANK	CLUSTER ID	SCORE	RNA	SCORE
1	Cluster 91401 ¹	0.45	MYB	0.53
2	Cluster 62518 ²	0.28	THSD4	0.24
3	Cluster 20099 ³	0.28	CT62	0.22

Table 2: Nestedness scores of the network. The NODF (BINMATNEST) score goes from 0 (form 100) for non nested network to 1 (to 0) for highly nested network. In the case of BINMATNEST temperature we have used 100 different random matrices and p-values for three null models are $p_1 = 0.01$, $p_2 = 0.04$ and $p_3 = 0.06$.

NODF	Temperature
0.3	4.4

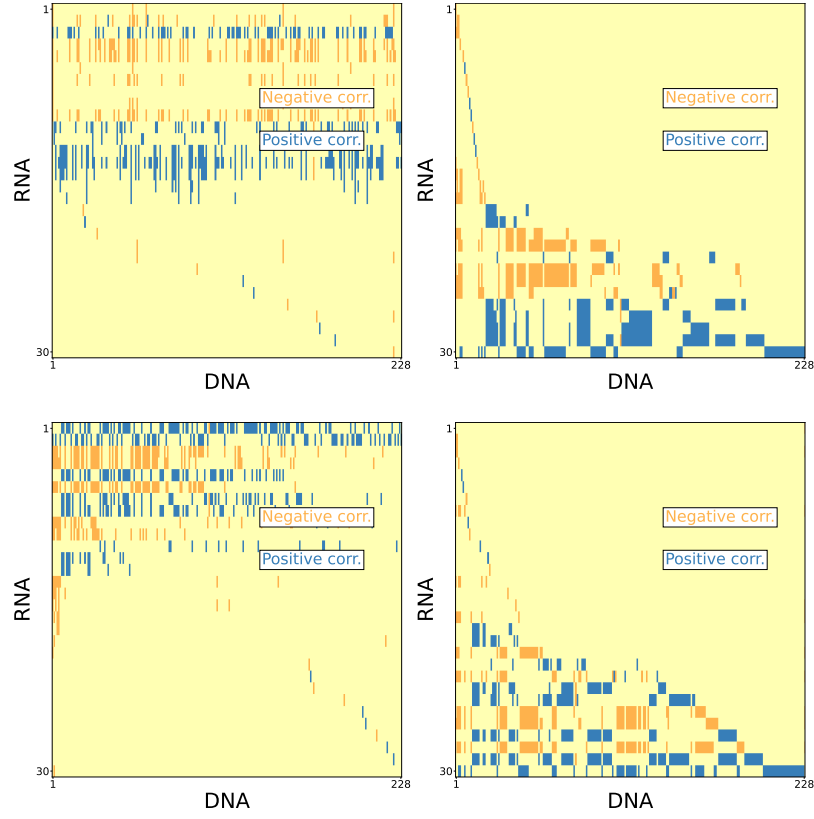


Figure 1: Binary bi-adjacency matrix representing the interactions between Clustered DNAs and RNAs involved in ESR1 synthesis and regulation. Matrix entries related to positive correlations are in blue, and negative in orange, finally empty cells are in bright yellow. Raw matrix (top left), Fitness and Complexity based reorganized matrix (top right), Degree based reorganized matrix (bottom left) and matrix with cols and rows reorganized from BINMATNEST algorithm. Here we have considered all Pearson correlation such that $|\rho| \geq 0.7$.