

Reconstruction of Gene Regulatory Networks using Topological Priors

*A thesis submitted in partial fulfillment
of the requirements for the degree of*

BACHELOR OF TECHNOLOGY

in

Electrical Engineering

by

Deepali Jain

Entry No. 2012EE10082

Suchakra Sah

Entry No. 2012EE10484

Under the guidance of

Sumeet Agarwal



**Department of Electrical Engineering,
Indian Institute of Technology Delhi.**

July-Nov 2015

Certificate

This is to certify that the thesis titled RECONSTRUCTION OF GENE REGULATORY NETWORKS USING TOPOLOGICAL PRIORS being submitted by DEEPALI JAIN and SUCHAKRA SAH for the award of B.Tech degree in Electrical Engineering is a record of bonafide work carried out by them under my guidance and supervision at the Department of Electrical Engineering. The work presented in this thesis has not been submitted elsewhere either in part or full, for the award of any other degree or diploma.

Prof. SUMEET AGARWAL
Department of Electrical Engineering
Indian Institute of Technology, Delhi

Abstract and Keywords

0.1 Abstract

This B.Tech project thesis explains how the Gene Regulatory Networks (GRNs) can be reconstructed from the corresponding gene expression data by assuming a scale free prior structure. It also focuses on algorithms used for generation of such data and the consequent sampling techniques. It then compares these techniques on the basis of Area Under the Curve within the ROC plot. Finally, the important results have been discussed and suggestions have been made for future work.

0.2 Keywords

GRN (Gene Regulatory Networks)

Gene expression data

Scale free prior

Markov Chain Monte Carlo (MCMC) approach

Barabasi Albert algorithm

Metropolis-Hastings Sampling

ARACNE (Algorithm for Reconstruction of Accurate Cellular Networks)

CLR (Context Likelihood or Relatedness Network)

MRMR (Maximum Relevance Minimum Redundancy)

Acknowledgments

Foremost, we would like to express our sincere gratitude to our advisor *Prof. Sumeet Agarwal* for the continuous supervision and support of this B.Tech project. This would not have been possible without his patience, motivation and immense knowledge.

We would also want to thank *our commitee members*, Prof S.D. Roy, Prof. Santanu Chaudhary, Prof. Seshan Srirangarajan and everyone for their insightful questions and comments.

Besides this, we would like to thank our *families*: our parents and sisters for their constant love and belief in us.

Deepali Jain
Suchakra Sah

Contents

0.1	Abstract	2
0.2	Keywords	2
1	Introduction	1
1.1	Gene Regulatory Networks	1
1.2	Gene Expression Matrix	2
1.3	Motivation for Reconstruction	3
1.4	Scale Free Networks	5
1.5	Problem Formulation	5
2	Theory	7
2.1	Literature Survey	7
2.2	Barabasi Albert Algorithm	8
2.3	MCMC Implementation	10
2.3.1	Posterior Probability	10
2.3.2	Metropolis Hastings	11
2.4	ARACNE	12
2.5	CLR Technique	13
2.6	MRNET Technique	13
2.7	MRNETB Technique	14
2.8	Genetic Algorithm	14
3	Results and Conclusions	15
3.1	Gamma estimation	15
3.1.1	Discussion	17
3.2	Threshold v/s γ	17
3.2.1	Discussion	17

3.3	AUC plots	17
3.3.1	Discussion	19
3.4	Comparison with Genetic Algorithm	21
3.4.1	Discussion	21
4	Conclusion and Future Work	23
	Bibliography	25

List of Figures

1.1	Typical GRN	2
1.2	Regulatory Interaction	2
1.3	Typical $p \times n$ gene expression matrix	3
1.4	Hubs and bottlenecks in genes	4
1.5	Scale Free Prior	6
2.1	Preferential Attachment <i>Source</i> : [14]	9
2.2	Metropolis Hastings	12
3.1	Synthetic dataset, Estimated $\gamma = 2.42$	16
3.2	Real dataset, Estimated $\gamma = 2.23$	16
3.3	Gamma v/s No. of Edges - 50 genes	18
3.4	Gamma v/s No. of Edges - 100 genes	18
3.5	ROC 50 genes, SN v/s 1-SP	20
3.6	ROC 100 genes, SN v/s 1-SP	20
3.7	ROC 30 genes, SN v/s 1-SP	22
3.8	ROC 40 genes, SN v/s 1-SP	22

List of Tables

3.1	γ Values	15
3.2	AUC values	19
3.3	AUC values, GA	21

Chapter 1

Introduction

With advancements in technology and wide interest in molecular biology, analysis of genomic data has continuously attracted scientists and researchers. The computations carried out has found use in diverse fields such as medicine, biology, therapeutics et al. In this chapter, we talk more about the problem statement and the potential applications of the same.

1.1 Gene Regulatory Networks

Gene Regulatory Networks, often referred as GRNs are a collection of genes in a cell which continuously interact with each other directly or indirectly (for example through their protein products) to ensure cell's function, fitness and metabolism.[3]

GRNs are usually represented as graphs where each node represents a gene and each edge represents an interaction between the two genes or nodes it connects. These regulations between the genes may be promotory or inhibitory in nature.

Fig 1.1 represents one such GRN. Note that the figure shows core regulators (pink) and their protein-interaction partners (yellow). Arrows from dashed ellipses indicate that the targeted nodes are regulated by all of the regulators present inside the ellipse. Some regulators appear multiple times in the network to reduce the number of intersecting arrows.

Fig1.2 further explains what these arrows mean.

Identification and study of these interactions in a GRN hep us to predict molecular pathays and regulatory relationships.

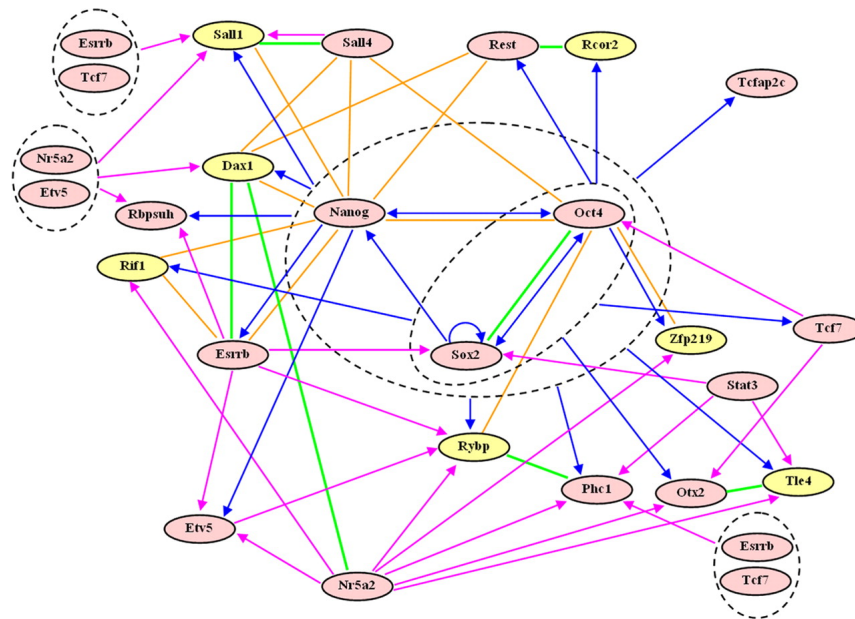


Figure 1.1: Source : [13]

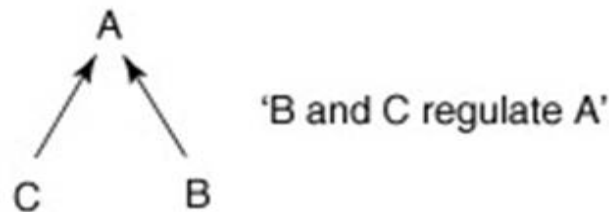


Figure 1.2: Source: [13]

1.2 Gene Expression Matrix

A cell contains thousands of genes each of which may or may not be activated at some given time instance i.e. they can be expressed in varying amounts within the cell depending on the task they are performing. Some genes expression level may be highly robust while for others genes, these levels may be very low.

These expression level of genes in a cell are measured by attaching large number of microscopic DNA spots on the surface under controlled experimental

conditions. The DNA sequence measured by each spot is called a Probe and the activation level of these Probes further indicate the expression level of genes. [6]

These levels when captured at different experimental conditions, say n samples for a set of p genes generate the $p \times n$ dimension gene expression matrix, as shown in Fig 1.3. This in turn can be used to reconstruct or learn the original gene network within the cell by reverse engineering.

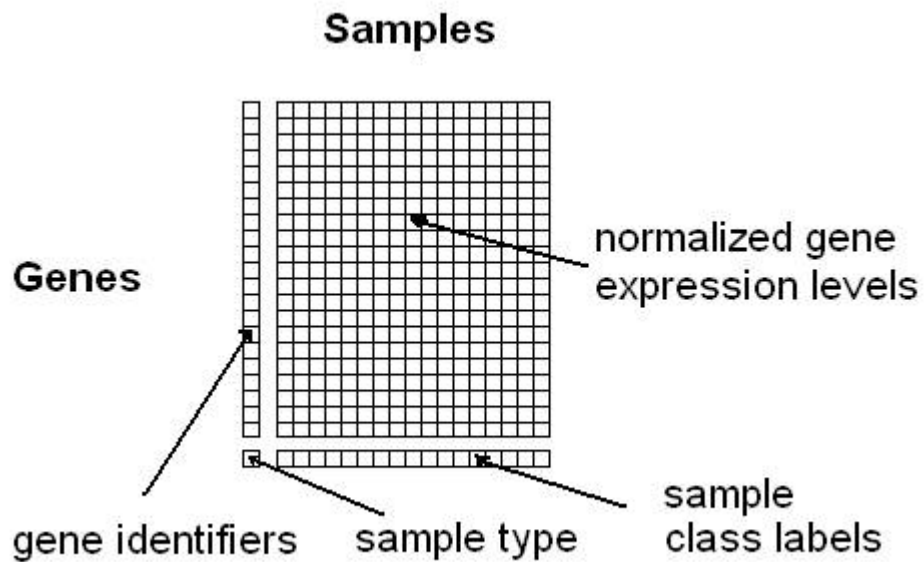


Figure 1.3: Source : [14]

1.3 Motivation for Reconstruction

The Gene Expression Matrix can be computed experimentally. From this matrix, we try to reconstruct the gene regulatory network i.e. a set of vertices V and edges E using various machine learning methods. Once the GRN has been estimated, it can serve the following purposes after statistical inferences[7] :

1. **Blueprint of Molecular Interactions :**

The GRNs serve as a 'map' or 'blueprint' for molecular interactions

and can give novel biological inferences about the same. Since the predicted edges in a GRNs represent actual physical binding between the genes, these networks contain valuable biological information. There are roughly 20,000 genes in humans and with the help of GRNs, certain hubs and bottlenecks i.e. genes which actively influence other genes can be indentified to narrow down the potential interactions for effective experimentation techniques. Such hub and bottlenecks are best shown in Fig 1.4.

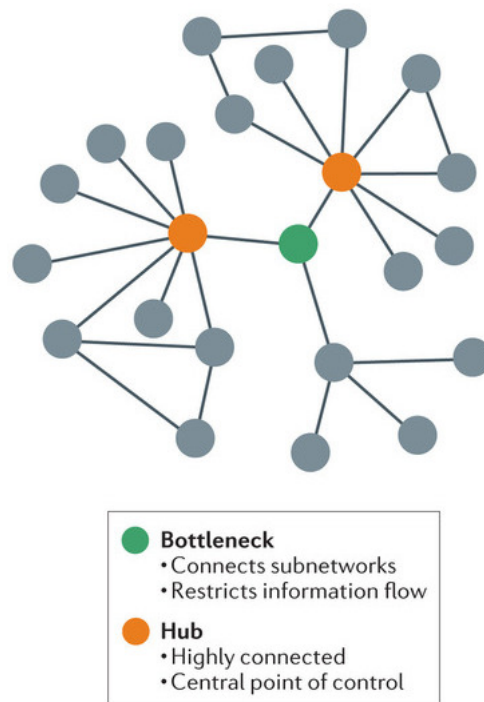


Figure 1.4: Source : [1]

2. Medicinal Uses :

Certain diseases such as cancer can be studied better by identifying the molecular pathways rather than the individual genes. Biomarkers based on individual genes neglect these pathways. These are however beautifully captured by GRNs allowing targeted medicines to be made. For example, gene expression data taken from a breast cancer tumor

study, shows that scale-free structure prior recovers hubs, including the previously unknown hub SLC39A6, which is a zinc transporter that is responsible for the spread of breast cancer to the lymph nodes.

3. Comparative Analysis :

When more and more GRNs corresponding to different disease conditions become available, we can statistically compare these graphs using techniques such as Graph Edit Distance. This is basically a measure of how similar two given graphs are. This will allow better study of diseases and infer the possible harms one condition can have on another.

1.4 Scale Free Networks

A scale free network is simply a network whose degree distribution follows a power law model i.e. there are large number of nodes with lesser links and few nodes with large number of links (acting as hubs) as shown in Fig 1.5. This is quite different from a random network where the degree distribution of a node is binomially distributed. Mathematically, the probability for a node in scale free network to have k edges, $p(k)$ can be stated as:

$$P(k) = k^{-\gamma} \quad (1.1)$$

where, γ lies in the range $[2, 3]$ for usual scale free GRNs.

1.5 Problem Formulation

There are limitations of the expression matrix. One, the number of independent experimental conditions have an upper bound and other, there is some noise inherent in the matrix. Hence, to increase the accuracy of the entire reconstruction process, we utilize certain previously established biological constraints or structural properties called priors along with the expression data. One property of these networks is being Scale Free.

In this thesis, we assume a Scale Free prior on the gene expression data and

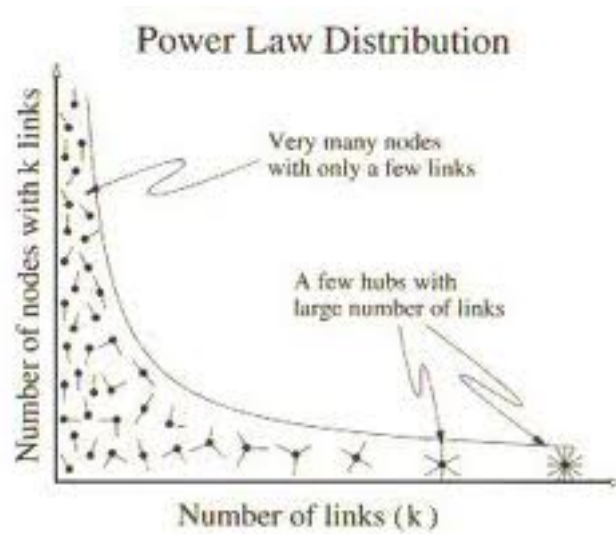


Figure 1.5: Source : [14]

then compare the results obtained with various reconstruction approaches such as Monte Carlo Markov Chain, ARACNE, CLR etc.

We then compare our work with Genetic Algorithm Approach implemented by Abdul [14].

Chapter 2

Theory

In this chapter, we talk about the main two publications that we referred to during the project. Then we explain the MCMC technique followed by talking about other reconstruction techniques against which we benchmark the MCMC approach.

2.1 Literature Survey

Since analysis and reconstruction of genomic data is now one of the latest research interests among scientists, there were many interesting papers and thesis available online. Though all of them have been mentioned in the *References* section, there are two papers from where our thesis draws its major inspiration. These are :

1. **A Scale-Free Structure Prior for Graphical Models with Applications in Functional Genomics** :[15] *Paul Sheridan , Takeshi Kamimura, Hidetoshi Shimodaira ; Nov 2010*

This research paper gives insights about the Markov Chain Monte Carlo Approach that has been followed by us. The static model which is a stochastic model depending on network parameter γ used for reconstruction has been proposed by Sheridan et al.

One of the other most important conclusions that has been drawn from this publication is how a scale free prior is better than a uniform prior assumption. A uniform prior is the one where the probability of a node having degree k in a p - node big network is binomially distributed. $P(k)$ is given as :

$$P(k) = \binom{p}{k} \beta^k (1 - \beta)^{p-k} \quad (2.1)$$

This implies that the actual degree in random network is centralized i.e. closer to the average node degree. Now when a random prior is used on a originally scale free network, the reconstruction obtained isn't as good as that obtained by assuming the scale free prior.

However the vice versa is not true. When we assume a scale free prior on a originally random network, the reconstruction is at par with a random prior assumption.

The real genomic dataset has also been provided with this publication and has been used by us to validate the MCMC approach by γ estimation.

2. Scale Free Prior in Gene Regulatory Network Reconstruction[14]

Abdul Hadi Shakir, IIT Delhi; June 2015

The M.Tech Project thesis by Abdul Hadi Shakir, a student of Computer Science department at IIT Delhi was another major source of knowledge and ideas throughout the making of our project. The performance metric used by us to compare various reconstruction techniques - the Area Under Curve in a ROC plot draws inspiration from Abdul's work. Abdul implemented the Genetic Algorithm technique for reconstruction and compared it with some state of the art methods available. He then tested this algorithm for synthetic data assuming uniform prior at first and then a scale free prior.

His thesis also concludes that scale free prior was able to give a better reconstruction of the network as compared to a uniform prior when AUC is used as a metric.

2.2 Barabasi Albert Algorithm

This algorithm is used for generation of scale free dataset for experimentation. Many common real networks are scale free such as World Wide Web, citation network and some social network. These networks share two popular mechanisms : Growth and Preferential Attachment. [2]

1. **GROWTH** : The network starts with an initial m_o nodes and with time, the network keeps on expanding in size i.e. more nodes are added to it.
2. **PREFERENTIAL ATTACHMENT** : When a new node enters the network, the probability of it being attached to a heavily connected node (node with a higher degree) is more than the probability of it being attached to a sparsely connected node, see Fig2.1. This basically implies that *Rich gets richer and poor get poorer*. It is an example of positive feedback where initially random differences in node degree are reinforced or magnified with time.

Mathematically, the probability of the new node being connected to an existing node i is p_i give as :

$$p_i = \frac{k_i}{\sum_j k_j} \quad (2.2)$$

where k_i is degree of node i . The degree distribution resulting from the BA is generally scale free with γ approximately 3.

SYSGENSIM uses this technique for dataset generation. It calculates the total links in the network by multiplying average node degree with total number of nodes. Then these links are distributed across the network to obtain scale free property.[4]

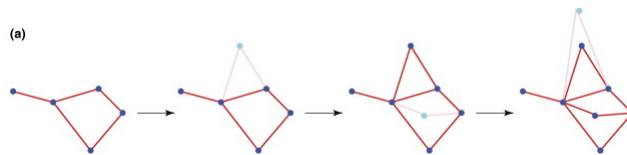


Figure 2.1: Preferential Attachment
Source : [14]

2.3 MCMC Implementation via Metropolis-Hastings Sampling

For a Bayesian setting,

$$P(G|D) \propto P(D|G) \times \pi(G) \quad (2.3)$$

where, G is the Gene graph, D is the dataset (gene expression dataset) and prior is scale free. Hence,

$$\text{Posterior Probability} \propto \text{Maximum Likelihood} \times \text{Prior}$$

The prior used by us is the scale free prior.

2.3.1 Posterior Probability

Let, σ be a **node labeling** or a permutation of integers $(1, 2, ..p)$ where p as discussed in Chapter 1 is the number of genes in the network. When each node is represented by integer σ_i , then $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_p)$.

We further proceed by **attaching a weight** w_i to each node v_1, v_2, \dots, v_p of the graph

$$w_i = \frac{\sigma_i^{-\mu}}{\sum_{l=1}^p \sigma_l^{-\mu}} \quad (2.4)$$

where μ is a tunable paramter in $(0,1)$. An important relationship is that between μ and γ (power law distribution parameter) given by :

$$\gamma = 1 + \frac{1}{\mu} \quad (2.5)$$

$$(2 < \gamma < 3)$$

When a GRN has to be generated, nodes v_i and v_j are selected with respective probability of w_i and w_j . These nodes are connected except when they are already connected or $i = j$.

The forementioned step is repeated pXK time where $p^{-\mu} \ll K \ll p^{1-\mu}$. K

is a parameter controlling the average number of edge.

The new model parameter is now $\theta = (\mu, K)$.

Further assuming that two nodes are selected independent of each other with probability = w_i , the **probability of two nodes i, j not being connected**, $P_{i,j}$ is calculated as :

$$P_{i,j} = (1 - 2w_i w_j) \sim (e^{-2pKw_i w_j}) \quad (2.6)$$

Now the **posterior probability for the generated network** is given by product of probabilities of edges present in graph and those not present in the graph. [14]

$$P(G|\theta, \sigma) = \prod_{\{v_i, v_j\} \in E} (1 - P_{i,j}) \prod_{\{v_i, v_j\} \notin E} (P_{i,j})$$

$$P(G|\theta, \sigma) = \prod_{\{v_i, v_j\} \in E} (1 - e^{-2pKw_i w_j}) \prod_{\{v_i, v_j\} \notin E} (e^{-2pKw_i w_j})$$

Hence, given the graph parameters, we can find out the posterior probability of a graph and compare it.

2.3.2 Metropolis Hastings

One problem that is faced is the number of possible networks for a given set of genes p , increases exponentially with size of p . We cannot check our results on every possible graph. Hence we need to use some sampling technique to sequentially arrive at the desired graph.

1. The Metropolis-Hastings algorithm works by generating a sequence of network parameters such that, as more values are produced, the distribution moves closer to the desired distribution.
2. At each step, it either adds or removes a edge from the graph and then accepts the updated graph if its posterior probability given by Eqn 2.5 is higher than an acceptance threshold.

In Fig 2.2, if the posterior probability is within the green range, the graph is updated else the proposed changes are rejected and some other edge is added or deleted from it.

3. The hyperparameter θ_m is updated by choosing a value with uniform probability distribution from $(\theta_m - \epsilon, \theta_m + \epsilon)$ for a step size of ϵ . This update is rejected if the parameter lies outside the range $(\theta_{min}, \theta_{max})$.

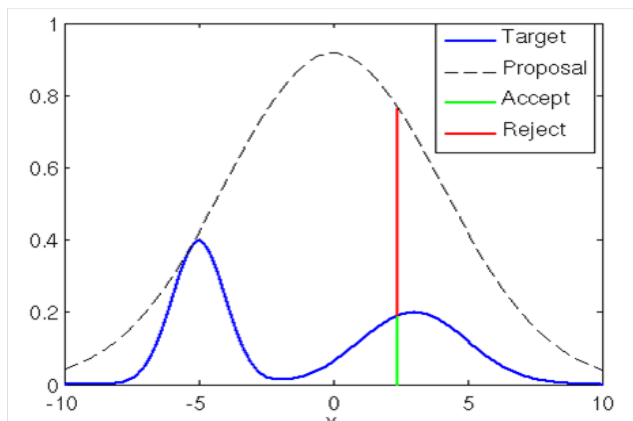


Figure 2.2: Source : [11]

2.4 ARACNE

ARACNE (Algorithm for the Reconstruction of Accurate Cellular Networks) is a state-of-the-art method for network reconstruction by reverse engineering. It defines an edge as an "irreducible statistical dependency between gene expression profiles that cannot be explained as a result of other statistical dependencies in the network".[10]

It further assumes that these dependencies can be inferred by studying the **pairwise statistical information** and higher order analysis are not required. This allows the algorithm to return a subset of the true regulatory interactions with less number of false positives.

The **mutual information for two genes** i, j is then calculated as $MI_{i,j}$ using :

$$MI_{i,j} = S_i + S_j - S_{i,j} \quad (2.7)$$

where S is defined as the entropy, a measure of randomness.

It then considers each triplet of nodes and removes the least significant edge in each triplet. In other words, when the $MI_{i,j}$ is above a threshold MI_o , an edge is created between gene i and j .

It can be implemented using 'minet' R package.[12]

2.5 CLR Technique

CLR (Context Likelihood or Relatedness Network)[9] instead of considering only the mutual information between a pair of genes i.e. $MI_{i,j}$, it also considers a scoring measure defined as $\sqrt{z_i^2 + z_j^2}$ where z_i for gene i can be obtained as

$$z_i = \max\left(0, \frac{MI(x_i; x_j) - \mu_i}{\sigma_i}\right) \quad (2.8)$$

and μ_i and σ_i are the mean and standard deviation resulting from computing all mutual information values involving the target as well as the regulator.

The rationale behind this method is the fact that most of the individual MI values involving the target or regulator are usually sparse and insignificant.

2.6 MRNET Technique

MRNET approach basically repeats MRMR (Maximum Relevance Minimum Redundancy)[8] feature selection procedure for every dataset variable/edges. The MRMR approach works by updating a set of selected variables S by adding a variable X_k such that

Minimize Redundancy,

$$W_k = \frac{1}{|S|^2} \sum_{i,j \in S} MI_{i,j} \quad (2.9)$$

and Maximize Relevance,

$$V_k = \frac{1}{|S|} \sum_{j \in S} MI_{h,j} \quad (2.10)$$

where h is the target distribution. Hence we maximize $V - W$.

2.7 MRNETB Technique

It is very similar to MRNET in terms of maximizing relevance and minimizing redundancy, but it works as a backward elimination process. It starts with a huge set of features and step-by-step eliminates those which follow the MRMR criteria. An important advantage of this approach over MRNET is that the former efficiently preserves those set of features which when present together significantly increase the relevance or prediction accuracy. [8]

2.8 Genetic Algorithm

In Genetic Algorithm approach (GA) we start with an initial subset of the population. Successive generation populations are generated by sharing of information with previous population. This can either be done by cross-over or mutation [5]. If the fitness of this new network is more than the median fitness of previous population, the network is added to the next generation population. The algorithm stops once we reach a convergence criteria or after a fixed number of iterations. The convergence criteria is usually set as constant median fitness of the population over 10 iterations. [14]

Chapter 3

Results and Conclusions

Firstly we validate the accuracy of the MCMC approach by matching the experimental and actual γ of synthetic dataset. Then the γ for real dataset is estimated.

We then observe how γ changes with the number of edges in the generated graphs. Finally we compare the reconstruction techniques explained in previous chapter.

3.1 Gamma estimation

For a scale free network, we have :

$$P(k) \propto k^{-\gamma} \tag{3.1}$$

$$\log(P(k)) \propto -\gamma (\log(k))$$

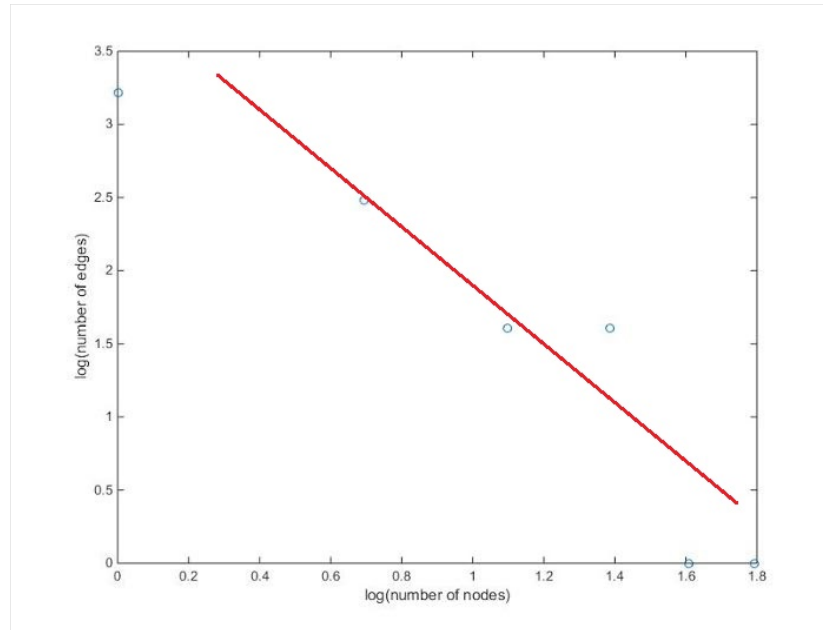
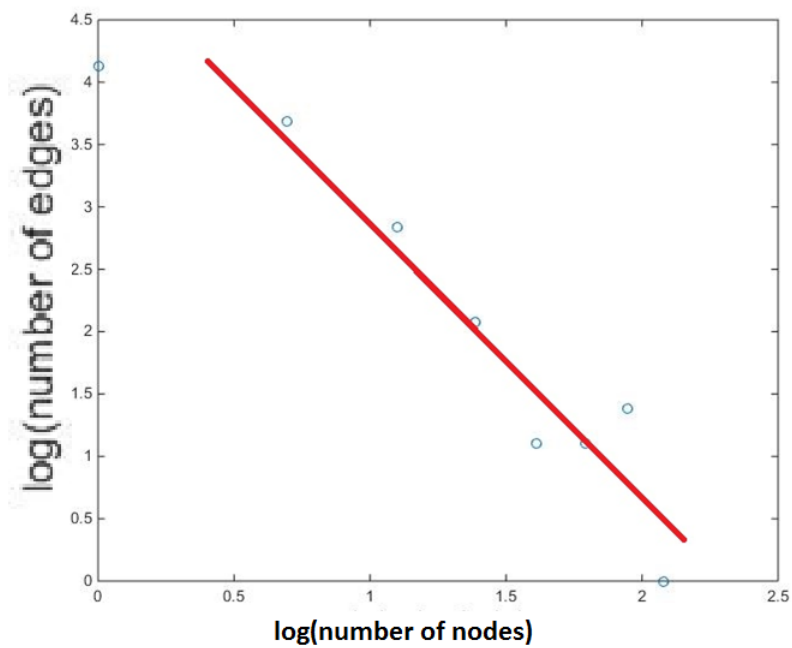
Hence, the graph of **log(probability of node degree = k) v/s log(k)** is **a straight line with a negative slope = $-\gamma$** .

Fig 3.1 and 3.2 show these straight line plots for synthetic dataset generated using SYSGENSIM and real dataset provided by Sheridian et al.

On comparing the γ values for these datasets, we obtain Table 3.1 results :

Table 3.1: γ Values

Dataset	γ Actual	γ Estimated
Real dataset	2.28(Sheridian et al)	2.23
Synthetic dataset	2.5	2.42
Synthetic dataset	2.7	2.61

Figure 3.1: Synthetic dataset, Estimated $\gamma = 2.42$ Figure 3.2: Real dataset, Estimated $\gamma = 2.23$

3.1.1 Discussion

We observe that the results are as one could expect. The γ estimated is almost 2.5 and the plot are straight lines with negative slope. This validates the MCMC approach used by us for scale free network reconstruction.

3.2 Threshold v/s γ

As discussed in previous Chapter for MCMC implementation, threshold is the value above which the graph posterior should lie after adding or deleting an edge for the graph to be updated.

We vary this threshold and plot the values of gamma obtained v/s the number of edges present in the resulting graph.

3.2.1 Discussion

It is observed in Fig3.3 and 3.4 that when the number of edges in resultant graph is close to edges in original graph, γ is closest to 2.5. This is the ideal value for γ in a scale free network. It is also observed that in all the cases, γ lies more or less in the range $[1, 5]$ and does not reach arbitrarily high values. Hence, even when uniform networks are estimated by assuming scale free prior, the results are at par by those obtained by assuming uniform prior itself as illustrated in the Literature Survey. The samples or experimental conditions used are 500 for each in the Fig3.3 and 3.4.

3.3 AUC plots

The MCMC, ARACNE, CLR, MRNET and MRNETB approaches for graph reconstruction are compared amongst each other on basis of their AUC value. This performance metric is defined as follows :

1. N_{TP} : Number of true positive edges in reconstructed graph

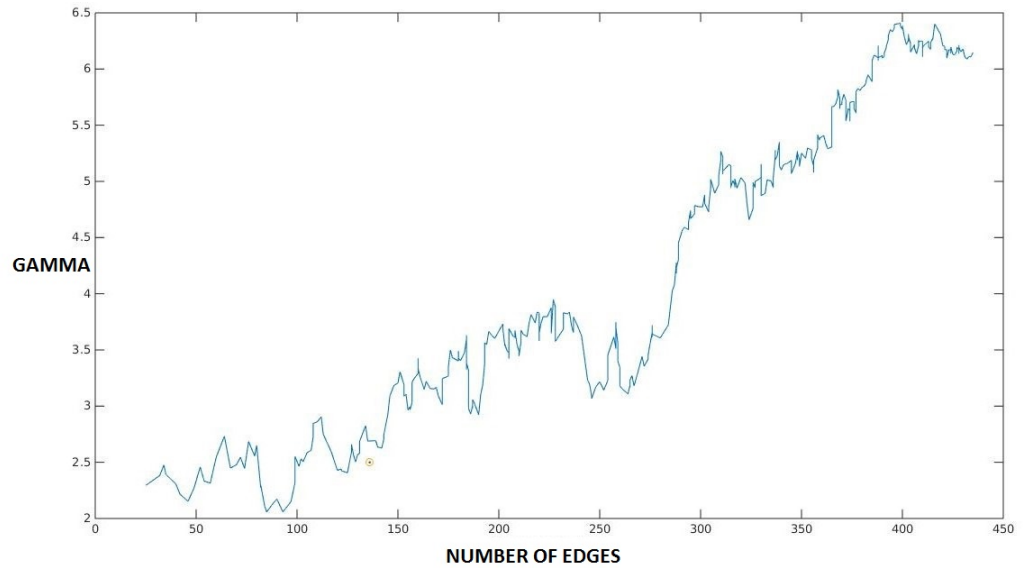


Figure 3.3: Gamma v/s No. of Edges - 50 genes

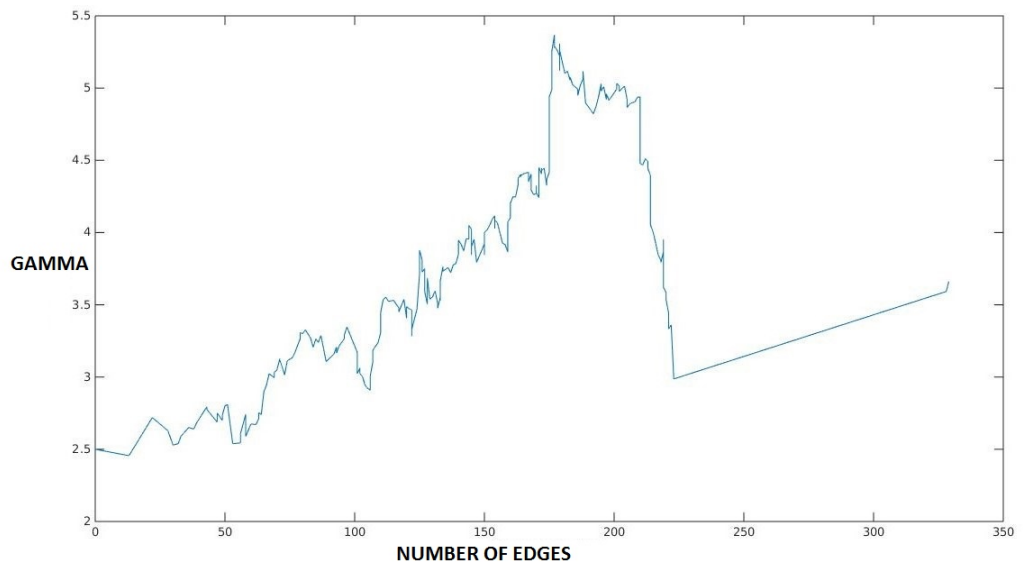


Figure 3.4: Gamma v/s No. of Edges - 100 genes

2. N_{TN} : Number of true negative edges in reconstructed graph
3. N_{FP} : Number of false positive edges in reconstructed graph
4. N_{FN} : Number of false negative edges in reconstructed graph

$$\text{Sensitivity, } SN = \frac{N_{TP}}{N_{TP} + N_{FN}}$$

$$\text{Specificity, } SP = \frac{N_{TN}}{N_{FP} + N_{TN}}$$

which means that SN is probability of detecting an edge given that it is present in the original graph while SP is the probability of not detecting the edge given that it is absent in the original graph.

Receiver Operator Characteristics (ROC) graph is a plot for SN against 1-SP. As we vary the threshold, we obtain different SN and SP values and the area under the ROC changes. More the AUC, better the reconstruction i.e. recovery of original network. The AUC should lie above the $x = y$ line otherwise the learning methods provide no insights.

Fig 3.5 and 3.6 show these AUC plots for varying node size and 500 samples or experimental conditions. These values are tabulated in Tab 3.2.

Table 3.2: AUC values

Nodes	MCMC	ARACNE	CLR	MRNET	MRNETB
50	0.744	0.921	0.8850	0.931	0.926
100	0.553	0.751	0.779	0.775	0.748

3.3.1 Discussion

We observe that ARACNE and MRNET are state-of-art techniques. MCMC approach has AUC values quite less than the techniques it is compared with. One of the possible reasons for its under-performance is that it does not give appropriate weightage to the fact that addition of two or more edges collectively may considerably increase the posterior probability of the graph. It adds or removes the edges from the graph one by one. On the contrary, this type of behaviour in edges forms the very basis of the MRNET and MRNETB techniques.

Underperformance of MCMC when compared to ARACNE can be attributed

to the fact that it does not take into account the pairwise mutual information between nodes. Also, the AUC values from ARACNE never reaches 1 since it rejects those pair of interactions with very low mutual information value. CLR again, extends this mutual information technique.

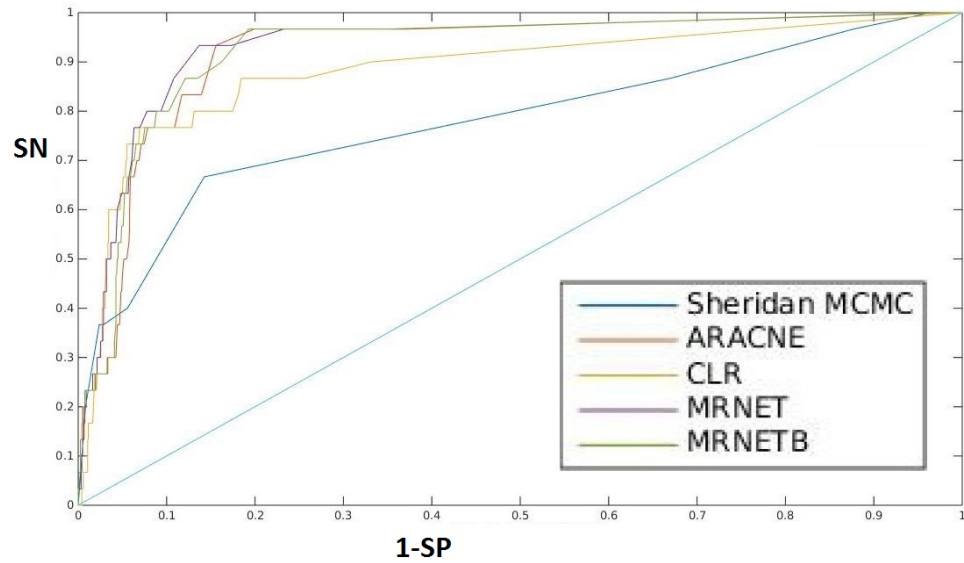


Figure 3.5: ROC 50 genes, SN v/s 1-SP

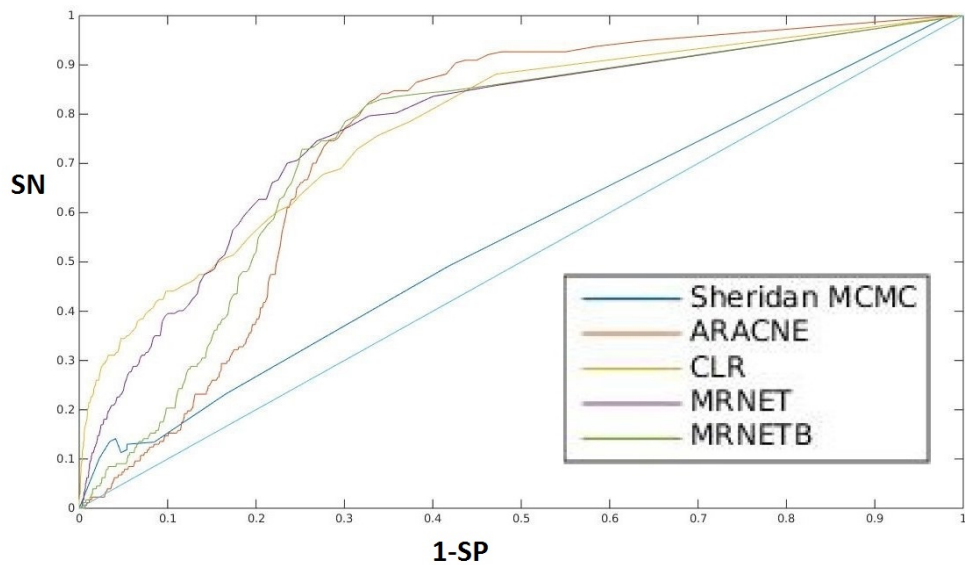


Figure 3.6: ROC 100 genes, SN v/s 1-SP

3.4 Comparison with Genetic Algorithm

Finally we compare these approach with Genetic Algorithm implemented by Abdul [14].

The genetic algorithm approach is implemented using DDEPN R package. The comparison for GA and MCMC is done for 30 and 40 genes. The performance metric used is again AUC value. We used this metric for various approaches on the dataset used by Abdul for his analysis.

Table 3.3: AUC values, GA

Nodes	MCMC	GA	ARCNE	CLR	MRNET	MRNETB
30	0.639	0.532	0.721	0.712	0.707	0.693
40	0.539	0.515	0.584	0.604	0.602	0.599

3.4.1 Discussion

We observe that GA and MCMC is not as good as techniques used in previous section - ARACNE, CLR, MRNET, MRNETB. The AUC values differ considerably. This further leads to the conclusion that it is important for reconstruction algorithms to take into account the collective action of two or more genes. Also, MCMC outperforms GA approach by a considerable amount. One of the possible reasons for this observation is that GA takes into account the fittest individuals and the overall fitness of the networks. This fitness measure favours sparse networks. [5]

The ROC are plotted as in Fig 3.7 onwards. The sample size is 500 and MCMC approach is used.

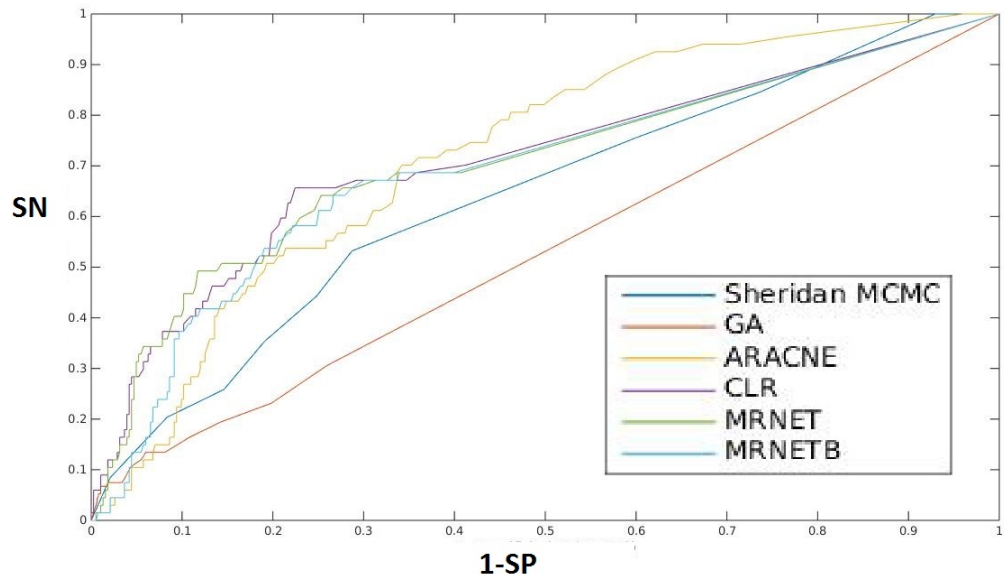


Figure 3.7: ROC 30 genes, SN v/s 1-SP

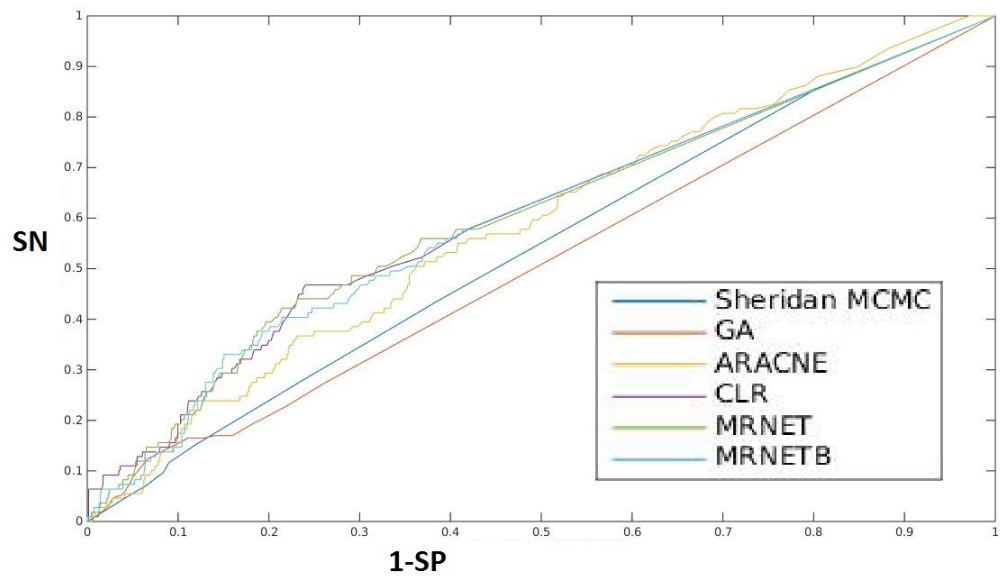


Figure 3.8: ROC 40 genes, SN v/s 1-SP

Chapter 4

Conclusion and Future Work

A conventional approach is to estimate the original gene network using the gene expression data by maximising the likelihood without incorporating the known priors. This may give us unrealistic networks. Hence, we are motivated to use some known biological constraints in our learning methods, one such constraint being the scale free property of networks. The Literature review clearly illustrates that how scale free prior perform better than other possible priors such as uniform priors.

Having understood this property, we discuss the scale free network generation algorithm (Barabasi Algorithm) and the sampling methods that can be used to reconstruct them. The basic idea behind other reconstruction methods prevalent such as ARACNE, CLR, MRNET, MRNETB are also highlighted in the theory chapter.

Finally, we observed that using the MCMC approach, we could obtain a straight line plot between the $\log(\text{probability of node degree being } k)$ and $\log(k)$. The slope of this line is $-\gamma$. We then showed that as we vary the threshold and hence the number of edges of the resultant estimated graph, the γ learnt changes. When the number of edges in the original and resultant graph are most similar, γ is close to 2.5. We also compared the reconstruction methods on AUC performance metric. We found that our MCMC approach is still far behind other methods like ARACNE, MRNET, MRNETB or CLR. This is mainly due to the fact that MCMC does not explicitly try to take into account the fact that the collection of two or more edges may significantly add to the posterior probability for the reconstructed graph.

This work on gene regulatory network reconstruction using topological priors can be extended in number of ways. We can look into priors like *network motifs* which are basically small structure providing a specific functionality to the cell and tend to repeat more often in the gene network. We can drive

techniques like *Maximum Relevance Minimum Redundancy* towards scale free priors i.e. tweak them so that they take into account this constraints observed in most genomic data.

Bibliography

- [1] Systems virology: host-directed approaches to viral pathogenesis and drug targeting, *nature reviews microbiology* 11, 455-466 2013.
- [2] Wikipedia page, barabasi algorithm - https://en.wikipedia.org/wiki/Barab%C3%A1si%E2%80%93Albert_model.
- [3] Wikipedia page, gene regulatory network - https://en.wikipedia.org/wiki/Gene_regulatory_network.
- [4] I. Hoeschele A. de la Fuente A. Pinna, N. Soranzo. Simulating system genetics data with sysgensim, *bioinformatics* 27(17), pp. 2459-2462, 2011.
- [5] Christian Bender, Silvia, Frauke Henjes, Stefan Wiemann, Ulrike Korf, and Tim Beissbarth. Inferring signalling networks from longitudinal data using sampling based approaches in the R-package 'ddepn'. *BMC Bioinformatics*, 12(1):291+, 2011.
- [6] Gary A. Churchill. Fundamentals of experimental design for cDNA microarrays. *Nat Genet*, 32 Suppl:490–495, December 2002.
- [7] Frank Emmert-Streib, Matthias Dehmer, and Benjamin Haibe-Kains. Gene regulatory networks and their applications: Understanding biological and medical problems in terms of networks. *Frontiers in Cell and Developmental Biology*, 2(38), 2014.
- [8] Howard Hughes Medical Institute Hanchuan Peng Janelia Farm Research Campus. Minimum redundancy and maximum relevance feature selection and its applications.
- [9] Aviv Madar, Alex Greenfield, Eric Vanden-Eijnden, and Richard Bonneau. Dream3: Network inference using dynamic context likelihood of relatedness and the inferelator. *PLoS ONE*, 5(3):e9803, 03 2010.
- [10] Adam A. Margolin, Ilya Nemenman, Katia Basso, Chris Wiggins, Gustavo Stolovitzky, Riccardo Dalla Favera, and Andrea Califano. Aracne:

-
- An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics*, 7(S-1), 2006.
- [11] Cinzia Pedrazzoli (D-MATH) Mattia Bergomi. Bayesian statistics: Computational aspects.
- [12] Gianluca Bontempi Patrick E. Meyer, Frederic Lafitte. Mutual information networks, version 3.29.0.
- [13] 2007-10-16 PNAS journal. A gene regulatory network in mouse embryonic stem cells.
- [14] Abdul Hadi Shakir. Scale-free prior in gene regulatory network reconstruction, 2013.
- [15] Paul Sheridan, Takeshi Kamimura, and Hidetoshi Shimodaira. A scale-free structure prior for graphical models with applications in functional genomics. *PLoS ONE*, 5(11):e13580, 11 2010.