

# Simple Bayesian Gene Network Learning in Populus Draught Transcriptome Data

**Amir Almasi Zadeh Yaghuti**

Nanjing Forestry University

**Ali Movahedi** (✉ [ali\\_movahedi@njfu.edu.cn](mailto:ali_movahedi@njfu.edu.cn))

Nanjing Forestry University <https://orcid.org/0000-0001-5062-504X>

**Hui Wei**

Nanjing Forestry University

**Weibo Sun**

Nanjing Forestry University

**Mohaddeseh Mousavi**

Nanjing Forestry University

**Qiang Zhuge**

Nanjing Forestry University

---

## Research Article

**Keywords:** Gene network, Populus, microarray data, NCBI

**Posted Date:** March 4th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-287640/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 Simple Bayesian Gene Network Learning in *Populus*  
2 Draught Transcriptome Data

3  
4

5 Amir Almasi Zadeh Yaghuti <sup>1,†</sup>, Ali Movahedi <sup>1,†</sup>, Hui Wei <sup>1</sup>, Weibo Sun <sup>1</sup>, Mohaddeseh Mousavi <sup>1</sup>,  
6 and Qiang Zhuge <sup>1,\*</sup>

7

8 <sup>1</sup> Co-Innovation Center for Sustainable Forestry in Southern China, Key Laboratory of Forest  
9 Genetics & Biotechnology, Ministry of Education, College of Biology and the Environment, Nanjing  
10 Forestry University, Nanjing 210037, China; amir\_20364@yahoo.com (AAZY);  
11 ali\_movahedi@njfu.edu.cn (A.M.); 15850682752@163.com (HW); czswb@njfu.edu.cn (WS);  
12 m.m2132@yahoo.com (MM)

13 \* Correspondence: qzhuge@njfu.edu.cn; Fax: +86-25-8542-8701

14

15 † These authors are contributed equally to the first author.

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34 **Abstract**

35 *Populus* is not only important for wood-based products, such as paper and timber, but also for  
36 metabolic-based production, for instance, bioethanol and biofuels. Constructing a sensibly  
37 functional gene interaction network is highly appealing to better understand system-level  
38 biological processes governing various *Populus* traits. Bayesian network learning provides an  
39 elegant and compact statistical approach for modeling causal gene-gene relationships in  
40 microarray data. Therefore, it could come with the illumination of functional molecular playing  
41 in Biology Systems. In this study, different forms of gene Bayesian networks were learned on  
42 *Populus* cellular transcriptome data. We addressed that Markov blankets, separating genes  
43 external to a regulatory Bayesian network from its internal genes, would likely be emerging at  
44 every possible gene regulatory Bayesian network level. The results have also shown that  
45 *PtpAffx.1257.4.S1\_a\_at,1.0* hypothetical protein is the most important in its possible regulatory  
46 program. This paper illustrates that the gene network regulatory inference is possible to  
47 encapsulate within a single BN model. Therefore, such a BN model can serve as a promising  
48 training tool for *Populus* gene expression data to better prepare future experimental scenarios.  
49

50 **Keywords:** Gene network; *Populus*; microarray data; NCBI

51

52

53

54

55

56

57

58

59

## 60 Introduction

61 The development of the theory of causal modeling roots back to the 1950s, which has left some  
62 scientific community controversies. The core of these debates and controversies is the Markov  
63 condition/ assumption, an assumption made in the Bayesian probability theory. Almost three  
64 decades later, causal modeling caught up with many researchers' attention, devoted to  
65 developing the theory of Bayesian networks(Gillies and Sudbury 2013). BN has been used to  
66 model chlordecone bioaccumulation in plants (Liber, Cornet, et al. 2020) to discover the best  
67 regulators of drought response (Lahiri, Venkatasubramani, et al. 2019) and to infer gene  
68 regulatory networks (Vignes, Vandel, et al. 2011). A Bayesian gene network consists of a  
69 digraph, which connecting regulatory genes to their targets, and elegantly encodes conditional  
70 independence between genes. Technically, BN is a combination of Bayesian theorem with the  
71 directed acyclic graph (DAG). The DAG decomposes the joint probability distribution. Given  
72 the DAG, the joint probability distribution of the expression of 6 nodes/genes factorizes as  
73 follows:  $P(A, S, E, O, R, T) = P(A)P(S)P(E | A, S)P(O | E)P(R | E)P(T | O, R)$

74 In general terms, this can be written as:  $p(x) = \prod_{i=1}^n p(x_i | \text{parents}(x_i))$ . In a BN, every  
75 gene must have a conditional probability table (CPT). For each gene, CPT indicates all the  
76 possible combinations of values of the parent genes. Each possible combination is called an  
77 Instantiation of the parent set. A BN structure is usually obtained from the data using score-  
78 based or constraint-based approaches (Koller and Friedman, 2009). Score-based algorithms  
79 maximize the BN likelihood, using Markov chain Monte Carlo (MCMC); to search the space  
80 of network structures, it operates on edge additions, deletions, or inversions. Score-based  
81 algorithms have been reported to fit well on simulated genetics and genomics data (Zhu,  
82 Wiener, et al. 2007) (Tasaki, Sauerwine, et al. 2015). However, we initially learned the  
83 network's undirected skeleton using repeated conditional independence tests in constraint-  
84 based approaches. Then by resolving directional constraints, each edge direction is assigned  
85 (v-structures and acyclicity) to the skeleton. Genetic mapping of multiple complex traits has  
86 been done(Scutari, Howell, et al. 2014). Due to its complexity, Bayesian gene network  
87 inference is feasible for systems of at most a few hundred genes or variables through  
88 conventional algorithms(Beckmann, Lin, et al. 2018, Wang, Audenaert, et al. 2019). The  
89 topology of the network of BN encodes conditional independence assertions. Therefore, we  
90 could logically find a set of possible gene regulators out of entire gene expression data. In this  
91 study, we try to find those genes in the *Populus* genome.

92 Drought is pivotal abiotic stress that affects plant development and poplar productivity  
93 (Hamanishi, Barchet, et al. 2015). Therefore, poplar breeding requires an understanding of the  
94 underpinning molecular regulatory machines controlling poplar resistance to drought stress.  
95 Molecular studies have revealed drought stress might take place in different plant tissues. For  
96 example, many transcriptomic studies, in poplar and also other plant genomes, have addressed  
97 transcriptional changes in roots (Molina, Rotter, et al. 2008, Cohen, Bogeat-Triboulot, et al.  
98 2010, Lorenz, Alba, et al. 2011, Stolf-Moreira, Lemos, et al. 2011, Dash, Yordanov, et al. 2018)  
99 due to drought stress over different genotypes (Cohen, Bogeat-Triboulot, et al. 2010, Stolf-  
100 Moreira, Lemos, et al. 2011, Cao, Jia, et al. 2014, Hamanishi, Barchet, et al. 2015, Jia, Li et al.  
101 2016) or at different levels of drought stress (Cohen, Bogeat-Triboulot, et al. 2010, Stolf-  
102 Moreira, Lemos et al. 2011). Gene network analysis on poplar root transcriptome identified a  
103 hierarchical-like gene network structure in which the highest hierarchical level 2,934 genes  
104 were affected by 9 super hubs (supergenomes) (Hamanishi, Barchet, et al. 2015).  
105 Here, we report microarray-based transcriptomic Bayesian gene network analysis of drought-  
106 resistant in black poplar genotypes measured in well-watered, moderate drought, severe  
107 drought, and post-drought re-watering conditions. Understanding the characteristics of poplar  
108 gene regulation of drought resistance, including various molecular interaction processes, shall  
109 elaborate functional knowledge of genes underlying this stress-induced phenomenon.

110

## 111 **Materials and Methods**

### 112 **Data**

113 Using the GEOquery package (Davis and Meltzer 2007), information directly from the GEO  
114 database with accession number GSE76322 was downloaded in this research. Initially, data  
115 dimensions were 61413×18. To run the Bayesian regulatory network, the probes with the  
116 highest variances were selected. Finally, the data dimension was reduced to 2210 × 18.

### 117 **Program**

118 We used several programs. The bnlearn package was used to infer the regulatory Bayesian  
119 network on our data. To determine the best network structure, the hill-climbing algorithm was  
120 used. The result of bnlearn was considered as a regulatory Bayesian network. The number of  
121 nodes, edges, Markov blanket size, neighborhood size, and the learned network's branching  
122 factor was calculated. In general, we assumed that the wider the Markov blanket size, the giant  
123 module network could be detected. With the bnlearn package help, the adjacency matrix with  
124 Cytoscape-based aMatReader software was imported to Cytoscape. Using the Network

125 Analyzer plugin, a comprehensive set of topological parameters such as number of nodes,  
126 edges, network diameter, radius and clustering coefficient, neighborhood connection, shortest  
127 path length, number of familiar neighbors, the degree distribution, the centripetal proximity  
128 parameter, etc. were calculated for our network.

129

## 130 Results and Discussion

131 In general, for high-dimensional genetics and genomics data, BN learning is challenging since  
132 the number of expected networks scales up exponentially with the number of genes. The  
133 computational cost of conventional BN inference could be a burden prohibitive (Wang,  
134 Audenaert, et al. 2019). Table 1 gives the general parameters of the observed probe-based  
135 Bayesian regulatory network. The number of probes (nodes), as can be seen in table 1, is 2210,  
136 which is higher than the number of directed edges (1000) that is logical in this context. This  
137 would indicate that some orphan probes in the network or some nodes are associated averagely  
138 high with other nodes.

139 One of the main parameters of detected BN is the Markov blanket (0.45).

140 Table 1: Estimation of structural Bayesian network parameters with Hill Climbing algorithm (MB:  
141 Markov Blanket, NS: Neighborhood Size, BF: Branching Factor, HC: Hill Climbing)

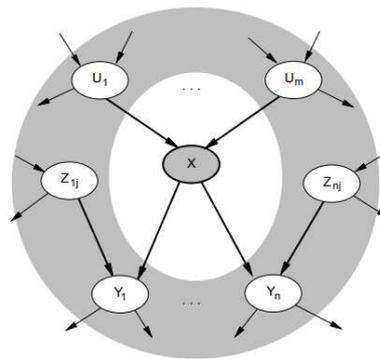
Parameters	values
No. of Nodes	2210
No. of Arcs (Edges)	1000
Undirected Edges	0
Directed Arcs	1000
MB	0.92
NS	0.90
BF	0.45
Penalization coefficient	1.445186
No. of Tests	4647736

142

143 To give a general biologically motivated expression of Markov blanket (MB), we could say that  
144 MB implicitly and explicitly is a minimum set of genes that would block all back-door paths  
145 that could be discovered from learned BN. The MB of a gene includes its parents, children, and  
146 co-parents. Parents in the MB separate the impact of specific disturbance on the outcome from  
147 the rest of the gene network. Therefore, this concept can be used to select a smaller set of  
148 relevant genes in high-dimensional problems. The MB is proven to be highly effective for  
149 feature reduction in high-dimensional problems, sometimes reducing the number of variables  
150 a thousandfold without any loss of accuracy (Aliferis, Tsamardinos, et al. 2003, Shen, Li et al.

151 2008, Fu and Desmarais 2010, Tan and Liu 2013). We put forward this idea that the MB  
 152 establishment is also instrumental in creating gene module networks. In other words, we could  
 153 say that number of MBs and their sizes in BN may reflect the sense of modularity in the  
 154 network. Parents of genes in the BN are the connections that would reflect causalities in the  
 155 gene network. By identifying the MBs in different sizes over all genes in the gene network, we  
 156 can quickly draw different gene sizes that cohesively function together, e.g., module network.

157



158

159 Figure 1: MB of a given postulated gene BN shows that MB of a node/probe/gene is the set  
 160 containing the node/gene's parents, children, and co-parents.

161 Figure 1 illustrates seven genes in a postulated gene regulatory BN. These seven genes  
 162 constitute a structure which is called MB. The  $U_1, U_m$  genes are parents of the  $X$  gene, and  
 163  $Y_1$  and  $Y_n$  genes are Children of the  $X$  gene in which the other two parents, e.g.,  $Z_{1j}$  and  $Z_{nj}$   
 164 genes they share paternity with the  $X$  gene. This structure demonstrates some cohesive acting  
 165 of genes, which we call them up as gene modules. Finding these structures in BN would be  
 166 biologically appealing. In our study, it was turned out that many MB are cohesively overlapped.  
 167 Therefore, a general picture of gene modules was not seen. We are expected to see some grape-  
 168 like structures, i.e., gene modules due to MB; such structures were not visually turned up on  
 169 the learned BN. The main finding of this study, e.g., *PtpAffx.1257.4.S1\_a\_at.1.0*, hardly  
 170 bordered as a specific set of probes as gene module. In other words, we could say that this  
 171 probe (and the gene it belongs to) might be involved in many other MB as well. This may  
 172 indicate the level of complication of droughtiness on the poplar genome. In other words, drought  
 173 stress has caused many different biological pathways actions. In every gene regulatory BN  
 174 network, we could find the total connectivity for each gene/probe, and this measure in BN can  
 175 be split up into out-degree and in-degree. Table 2 and figure 2 address some topological aspects  
 176 of learned BN. Putting this in a biological context, we could see that nodulizing BN has been

177 hardly achieved. One reason could be the fact that drought stress affected many different  
 178 functional parts of the poplar genome.

179 However, the distribution of gene out-degree has been much higher than gene in-degree. This  
 180 is biologically supported as out-degree probes/genes could be considered as a regulator, which  
 181 patently indicates that a small number of genes had many out-degrees. Some of these high-out-  
 182 degree genes are shown in table3. The results of this table can be used for further analysis, for  
 183 example, utilizing suitable gene set enrichment analysis. A general feature of many gene  
 184 networks is their nature of scale-free topologies addressing a gene degree distribution that can  
 185 be shown with a function of power-law, e.g.,  $P(k) = Ck^{-\alpha}$ , in which  $P(k)$  is the randomly  
 186 selected gene having degree  $k$  (or  $k$  connections),  $\alpha$  is the power-law exponent, and the  
 187 constant  $C$  is a Riemann's zeta function e.g.  $\sum_{n=1}^{\infty} \frac{1}{n^s}$  ( $s$  = complex variable and  $n$  = integer)  
 188 normalizing the power law probability distribution, e.g.  $\sum_{k=1}^{\infty} P(k) = 1$ . In scale-free gene  
 189 network topologies, most genes generally have relatively few interactions reflected as a lower  
 190 degree. In contrast, a small number of genes, so-called 'hubs' genes, have a higher degree.

191 Table 2"simple parameters in network Analyzer.

<b>Clustering coefficient</b>	00.00
<b>Number of nodes</b>	2211
<b>Connected components</b>	1
<b>Network diameter</b>	1
<b>Network radius</b>	1
<b>Shortest paths</b>	2210(0%)
<b>Characteristics path length</b>	1
<b>The average number of neighbors</b>	1.999
<b>Network density</b>	0.0
<b>Isolated nodes</b>	426
<b>number of self-loops</b>	0
<b>multi-edge node pairs</b>	0
<b>Analysis times (Sec)</b>	0.874

192

193

194

195

196

197

198

Table 3: The CytoHubba plugin chose the first 30 genes in the MCC method.

<b>PtpAffx.1257.4.S1_a_at,1.0</b>	<b>Hypothetical protein /// hypothetical protein</b>
<b>PtpAffx.1258.1.S1_s_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.1258.4.S1_s_at,1.0</b>	<b>Aktin9</b>
<b>PtpAffx.1259.1.S1_s_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.12595.1.A1_s_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.126167.1.S1_at,1.0</b>	<b>Branched-chain amino acid aminotransferase-like protein</b>
<b>PtpAffx.126235.1.S1_at,1.0</b>	<b>Auxin-responsive family protein</b>
<b>PtpAffx.12628.1.S1_at,1.0</b>	<b>Metal handling. Binding, chelation, and storage</b>
<b>PtpAffx.126599.1.A1_s_at,1.0</b>	<b>Arginine decarboxylase</b>
<b>PtpAffx.1263.1.A1_a_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.12675.3.A1_s_at,1.0</b>	<b>Unknown protein; predicted by gene finder</b>
<b>PtpAffx.12675.3.A1_a_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.12684.1.S1_at,1.0</b>	<b>pyrophosphate-dependent 6-phosphofructose-1-kinase, putative</b>
<b>PtpAffx.12676.1.S1_at,1.0</b>	<b>MAC/Perforin domain-containing protein</b>
<b>PtpAffx.1269.1.A1_x_at,1.0</b>	<b>Putative thioredoxin</b>
<b>PtpAffx.1269.1.A1_a_at,1.0</b>	<b>Putative thioredoxin</b>
<b>PtpAffx.12712.1.S1_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.127056.1.A1_at,1.0</b>	<b>Signal transducer activity</b>
<b>PtpAffx.12745.1.A1_at,1.0</b>	<b>Unknown protein</b>
<b>PtpAffx.1273.1.S1_at,1.0</b>	<b>Copper ion binding, laccase activity,</b>
<b>PtpAffx.127847.1.A1_at,1.0</b>	<b>DNA binding protein</b>
<b>PtpAffx.127644.1.A1_s_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.12800.1.S1_a_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.127878.1.A1_at,1.0</b>	<b>Unknown protein</b>
<b>PtpAffx.1286.2.S1_s_at,1.0</b>	<b>S-adenosylmethionine decarboxylase</b>
<b>PtpAffx.1286.1.S1_s_at,1.0</b>	<b>S-adenosylmethionine decarboxylase, putative</b>
<b>PtpAffx.12874.1.S1_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.1286.5.S1_s_at,1.0</b>	<b>Hypothetical protein</b>
<b>PtpAffx.128895.1.S1_at,1.0</b>	<b>Lysine-ketoglutarate reductase/saccharopine dehydrogenase</b>

200

201

202

203

204

205

206

207

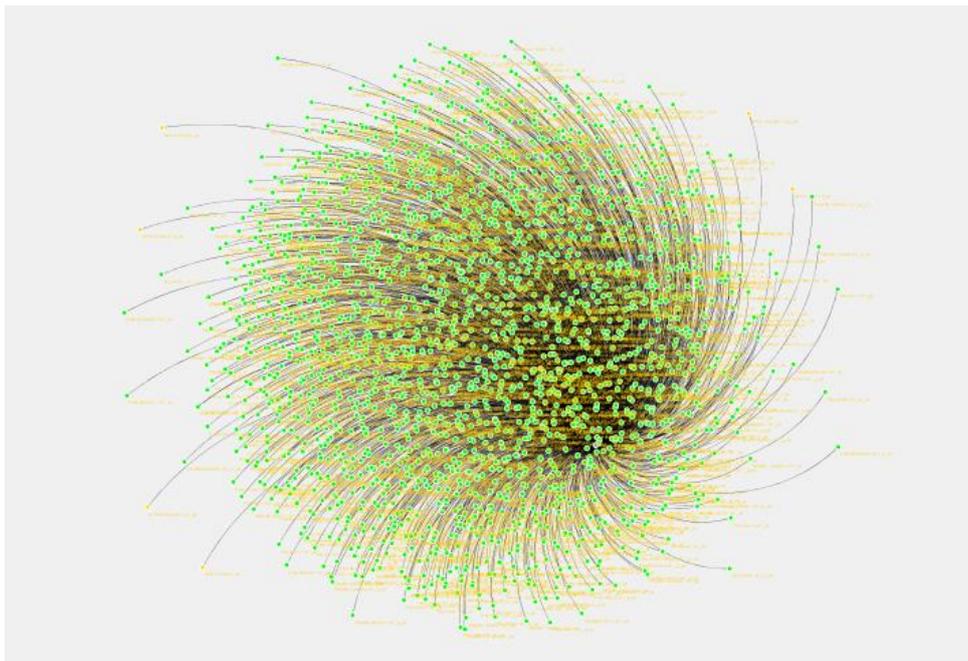
208

209

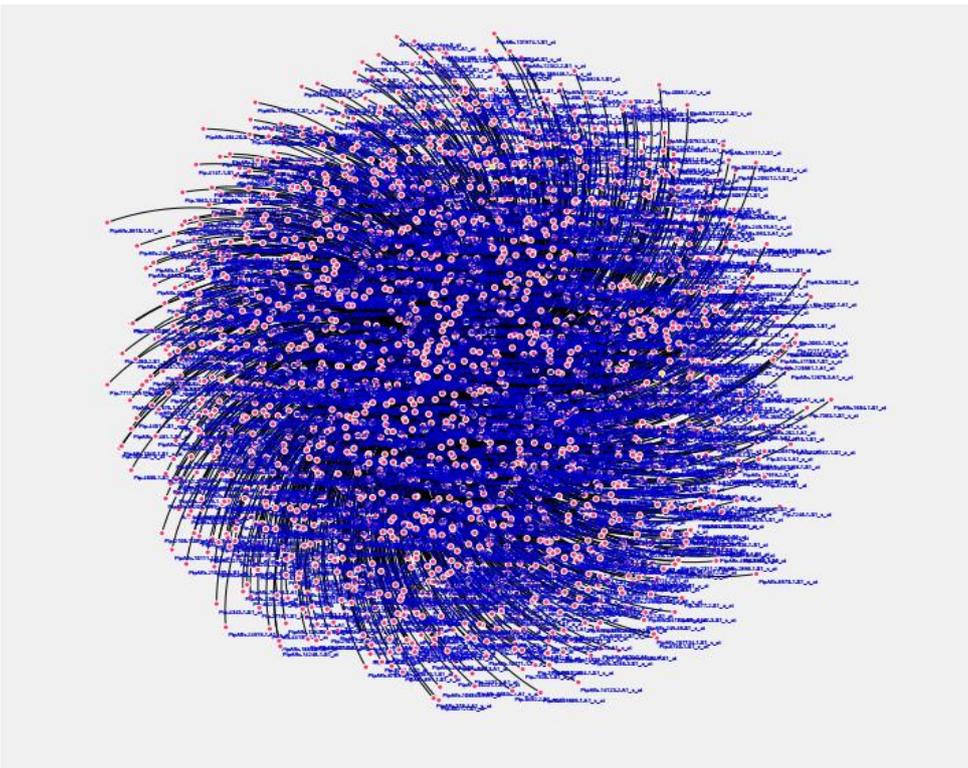
210

211

212 This may prohibit seeing clean-cut gene modules in studies similar to our study.



213



214

215 Figure 2: The snapshot of gene connectivity in learned BN in two graph forms.

216 To predict gene function in poplar gene expression data using DNA microarray datasets,  
217 private or publicly available databases are needed. Our method allows users to extract genes  
218 involved in biological processes, which could be valuable for understanding similar tree species  
219 mechanisms (Ogata, Suzuki, et al. 2009). The mode of genes in terms of being regulator or  
220 regulated is not reflected in this type of database. This could make the practical biological

221 application of such databases restrictive. We used public domain data in this study and to check  
222 up a new hypothesis in popups genome. This sort of data is being used to create a gene network  
223 (Cai, Li et al. 2014). In (Cai, Li et al. 2014), candidate *Populus* cell wall genes using a  
224 combination of a genome-wide *Populus* gene co-expression network (PGCN) and module  
225 detection and gene ontology (GO) enrichment analysis was performed to assign the functional  
226 category to these modules. However, our approach in this study was a probabilistic one instead  
227 of a Pearson-based gene correlation matrix. Our data did not find any topological properties of  
228 the network indicating scale-free and modular behavior. As it turned out, we were intensely  
229 relied on BN topological properties, e.g., Markov blanket, to address possible modules in our  
230 data. Figures 1 and 2 showed no sign of modules.

231 In the domain of gene expression high throughput data, probabilistic dependencies arise  
232 because genes are related in different ways (e.g., through common phenotype, logical  
233 connections, pathway connections, through (non-causal) physical laws, constrained of  
234 thermodynamic laws or structural boundary conditions). This may prohibit seeing clean-cut  
235 gene modules in studies like the current study. Figure 1 demonstrates that no sign of modules  
236 is trackable to pick one out of the gene regulatory network. Our method could likely help to  
237 identify and characterize cell wall-related genes in *Populus*. Using graph-based theory, (Zhang  
238 and Yin 2016) identified 14 probe-sets/genes related to the GO cellular component term of B  
239 plasmodesmata and many other genes for other tissues. However, the graph-based networks of  
240 the different tissues have shown different topological properties. They illustrate that genes in  
241 the root network were the most highly co-expressed. Whereas the leaf genes were the weakest  
242 co-expressed, and those in wood were in the intermediate. These topological network  
243 differences provide some unseen epigenetic mechanisms in different tissues. The genetic  
244 mechanism underlying the different topological properties among these tissues remains  
245 unknown and worthy of further investigation in future studies. In general, finding gene modules  
246 in *Populus* data has an almost long history (Grönlund, Bhalerao, et al. 2009, Liu, Ding et al.  
247 2016, Han, An et al. 2020). In most of them, the co-expression analyses have been the  
248 cornerstone in finding modules of co-expressed genes; gene modules display distinct  
249 expression profiles. However, addressing hub genes in this kind of study is a difficult task. For  
250 example, figure 2 shows how genes are cohesively related, and therefore, it is quite challenging  
251 to pick up and address some genes which actively responded to drought.

252  
253  
254

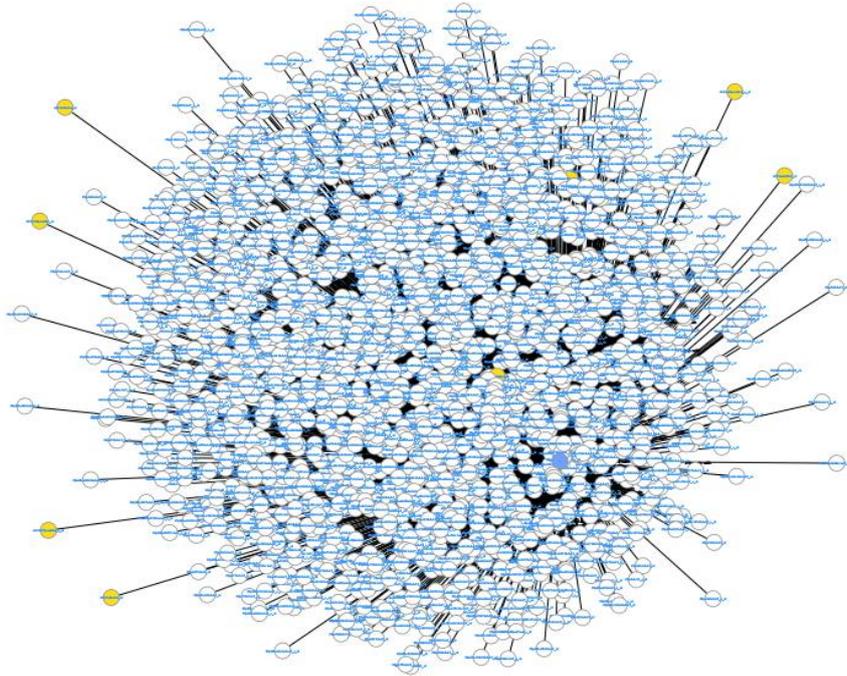
255  
256

Table 4. The identified DEC probe sets in leaf tissue involved in photosynthesis-related processes and light reactions

Biological process	Probe sets	Gene model	Description	TAIR gene model	Alias
Photosynthesis	PtpAffx.71066.5.A1_s_at	Potri.010G089400	Photosystem I 20kD family protein	AT1G03130	PSAD-2
Photosynthesis	Ptp.5240.1.S1_s_at	Potri.008G151600	Photosystem I 20kD family protein	AT1G03130	PSAD-2
Photosynthesis	Ptp.2148.1.S1_at	Potri.001G081500	Photosystem I chain III family protein	AT1G31330	PSAF
Photosynthesis	PtpAffx.249.164.A1_s_at	Potri.002G055700	Photosystem II oxygen-evolving complex protein two precursor	AT1G06680	OEE23, OEE2, PSBP-1, PSII-P
Photosynthesis	Ptp.1584.4.S1_x_at	Potri.005G239300	Chlorophyll a-b-binding protein 2	AT2G34430	LHB1B1, LHCB1
Photosynthesis	PtpAffx.689.1.S1_at	Potri.010G255500	Hypothetical protein	AT5G08410	FTRA2
Photosynthesis	Ptp.5386.1.S1_s_at	Potri.006G144000	Membrane protein	AT2G06520	PSBX
Photorespiration	PtpAffx.638.1.S1_s_at	Potri.010G045100	Alanine-2-oxoglutarate aminotransferase 1	AT1G70580	AOAT2, GGT2
Photorespiration	PtpAffx.638.1.S1_s_at	Potri.008G187400 Alanine-2-oxoglutarate aminotransferase 1	Photorespiration	PtpAffx.638.1.S1_s_at	Potri.008G187400 Alanine-2-oxoglutarate aminotransferase 1
Photorespiration	Ptp.1238.1.S1_at	Potri.009G081600	Malate dehydrogenase family protein	AT2G22780	PMDH1
Photosystem I	PtpAffx.7681.2.A1_at	Potri.009G065900	Hypothetical protein		
photosystem II	Ptp.6095.1.S1_at	Potri.006G088200	Hypothetical protein	AT5G02120	OHP
Light reaction	Ptp.8143.1.S1_s_at	Potri.011G099400	Cytochrome P450 family protein	AT3G14690	CYP72A15
Light reaction	Ptp.6095.1.S1_at	Potri.006G088200	Hypothetical protein	AT5G02120	OHP
Light reaction	Ptp.1584.4.S1_x_at	Potri.005G239300	Chlorophyll a-b-binding protein 2	AT2G34430	LHB1B1, LHCB1.4
Response to UV	Ptp.4373.1.S1_x_at	Potri.005G246100	Plastocyanin family protein	AT1G76100	PETE1

257

258



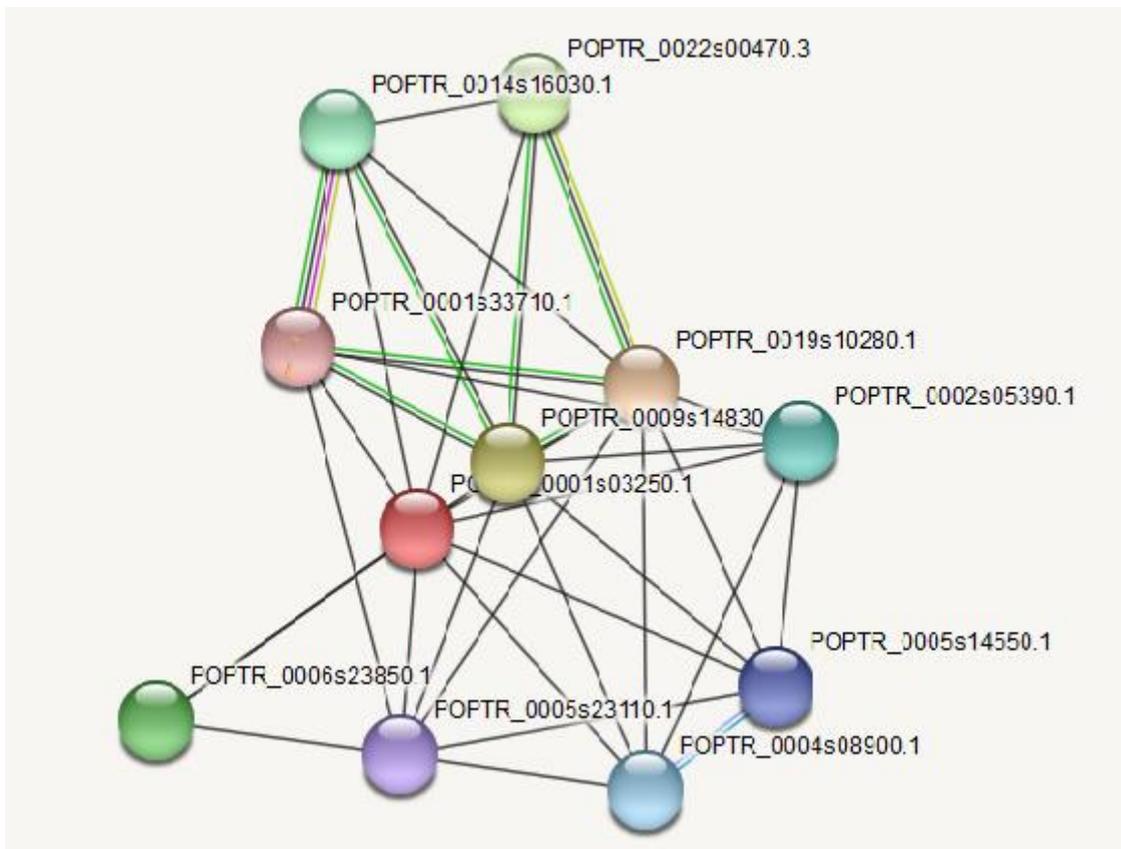
259

260 Figure 3: Highlighting top 10 genes with their neighbors using MCC algorithm

261 In the MCC algorithm, it is possible to highlight motivated biological molecules in the network.  
 262 Some yellow circle nodes in figure 3 display the top ten most connected hub genes. As it can  
 263 be seen, there is a close relationship between other genes and these top 10 genes, which in turn  
 264 makes it challenging to report some general gene module structures in our data. However, in  
 265 this study, the applied BN model should be regarded as a conjecture, which should be tested  
 266 using more and better posed high dimensional data.

267 Poplar has emerged as an ideal model system for studying woody plants. To better understand  
 268 the biological processes underlying various poplar traits, e.g., wood development, a  
 269 comprehensive functional gene interaction network is highly needed. To more effectively  
 270 screening down, we have used some highly connected gens in this study. Poplar Gene offers  
 271 comprehensive functional interactions and extensive poplar gene functional annotations.  
 272 Seminal research on the poplar genome shows that relatively small genome, quick growth  
 273 trend, and simple genome clonal manipulation are the interesting biological features that have  
 274 made poplar a long-lived forest tree model system (Liu, Ding, et al. 2016). Drought stress is  
 275 the leading cause of plant loss worldwide, and drought is one of the essential environmental  
 276 fluctuations that affect almost all plant species (Dash, Yordanov, et al. 2018). As a model,  
 277 *Populus* provides an opportunity to study the stress response in a perennial tree growing as a  
 278 commercial biomass product to produce carbon-neutral energy (Street, Skogström, et al.,  
 279 2006).

280 In this study, as can be seen in table 3, many genes turned out to be highly connected ones: hub  
281 genes. Moreover, their annotations reveal that they do not have comprehensive biological  
282 support (as indicated by hypothetical term). Therefore, it was not easy to see their possible  
283 protein interaction in the STRING database. Protein-protein interaction networks play an  
284 essential role in understanding the system level of cellular processes. These networks can filter  
285 and evaluate functional genomic data and create an intuitive platform for annotating proteins'  
286 structural, functional, and evolutionary properties. In this study, we took the aktin 3 genes  
287 (Table 3) and fed them to the STRING database, and the grid of protein interaction was  
288 extracted. STRING is a protein-protein interactions database (the known and predicted one)  
289 and a retrieval search tool. The results of drawing interactive networks in the STRING database  
290 are shown in figure 4. The result of the protein network was architecturally different from those  
291 obtained from the Bayesian gene network. In other words, some inconsistency and discrepancy  
292 in the results of the network learned by STRING were observed compared to the learned BN  
293 network, which shows that post-transcriptional modification may play crucial roles in  
294 regulating gene expression in the *Populus* genome.

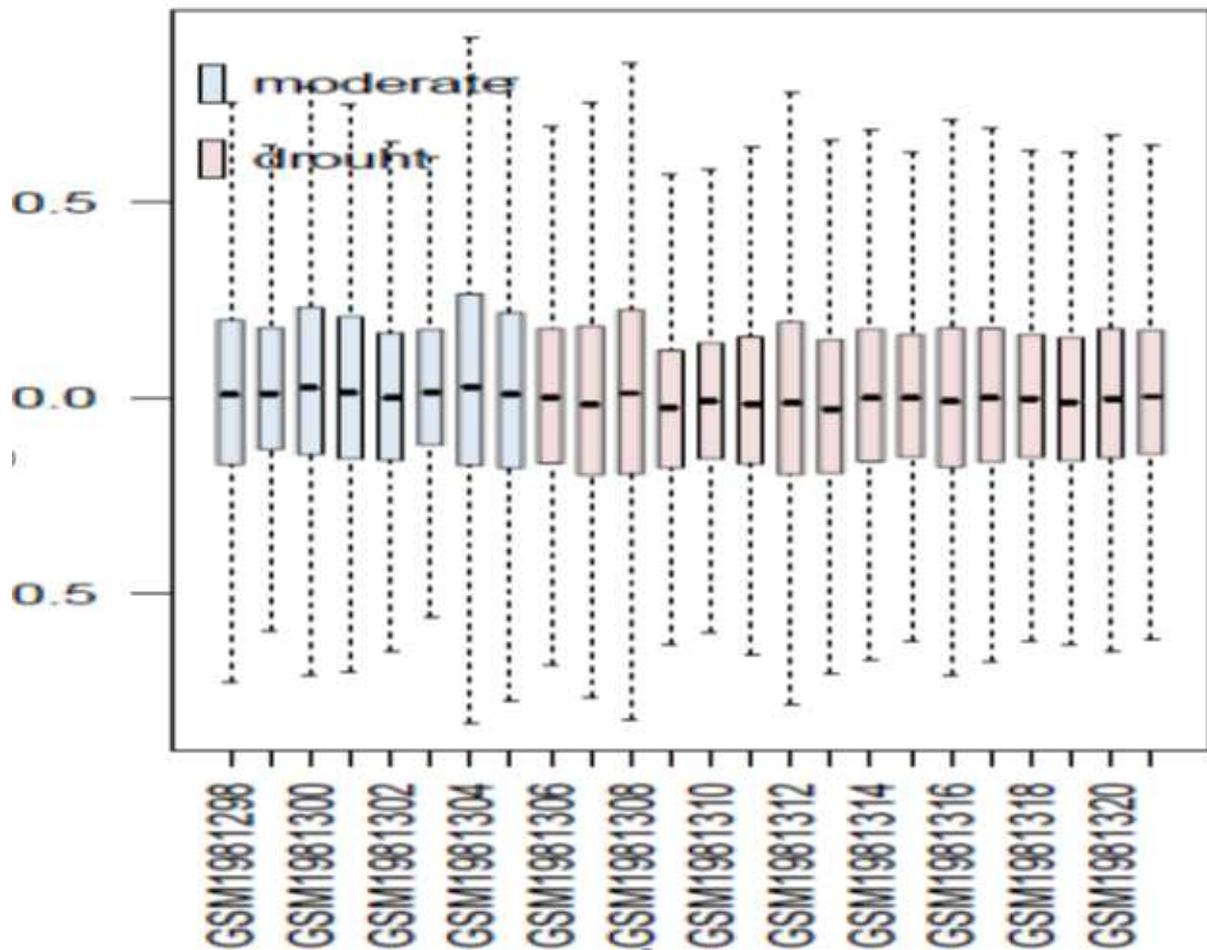


295 Figure 4: Network of protein-protein interactions by STRING database information.

296 In this study, we logically relied on *Populus* mRNA data; therefore, talking about the PPI  
297 network does not make that much sense. However, we believe that the whole body or tissue-  
298 specific protein-protein interactome (PPI) map could highlight drought tolerance in the  
299 *Populus*. Large-scale studies of PPIs accompanying the species-specific database support the  
300 interplay of cascading translational protein levels for drought-related stress adaptive  
301 mechanisms in *Populus*. Interestingly, many genes that became hub genes in this study were  
302 biologically involved in photosynthesis-related processes and light reactions. These are the key  
303 findings of this study. Most likely abiotic processes like drought act on the biological system  
304 level of *Populus*.

305 Boxplot is a visualization tool of microarray data. Each box stretches from the lower hinge -  
306 defined as the 25th percentile to the upper hinge - the 75th percentile- and the median is  
307 delineated as a line across the box. This figure has appeared in microarray experiment data at  
308 the top of other figures. However, we took some so-called downside-up approach, in a way that  
309 is shown at the bottom of the table to make it visually much more transparent than the  
310 treatments used in applied GSE76322, which did not impart much more diversity on the gene  
311 expression (figure 5). However, this sort of view should not mislead researchers, as in this  
312 study, many probes turned out to be differently expressed. We believe that BN has lots of  
313 potential in deciphering information on *Populus* microarray data, a matter that should be  
314 figured out in subsequent studies.

315



316

317 Figure 5: Boxplot of GSE764322 using in this study.

## 318 Conclusion

319 For the first time, a BN was learned out on *Populus* DNA microarray data. In this study, we  
 320 identified a gene network in the *Populus* to highlight candidate genes used as regulator genes.  
 321 We identified the top 30 hub genes, in which some of them had low if any, valid biological  
 322 information. We noted that the MB in BN of gene regulation could be assumed to narrow down  
 323 the whole gene network complexity. Statistically speaking, the existence of an MB means  
 324 external genes are conditionally independent of internal genes and vice versa. This is  
 325 biologically appealing. Despite extensive physiological and morphological descriptions of the  
 326 *Populus* response to drought, little work has been done to explain the differences in gene levels  
 327 and examine the similarity of the stress response between this perennial and the annual crop.  
 328 Because poplar genome sequences and poplar microarrays are now available, a bridge can be  
 329 made between quantitative trait locus mapping approaches, the candidate gene approach, and  
 330 transcription.

331 **Acknowledgment**

332 This study was supported by the National Science Foundation of China (No. 31570650) and  
333 the Priority Academic Program Development of Jiangsu Higher Education Institutions.

334 **References**

- 335 Aliferis, C. F., et al. (2003). HITON: a novel Markov Blanket algorithm for optimal variable  
336 selection. AMIA annual symposium proceedings, American Medical Informatics Association.
- 337 Beckmann, N. D., et al. (2018). "Multiscale causal network models of Alzheimer's disease  
338 identify VGF as a key regulator of disease." bioRxiv: 458430.
- 339 Cai, B., et al. (2014). "Systematic identification of cell-wall related genes in *Populus* based on  
340 analysis of functional modules in co-expression network." PloS one 9(4): e95176.
- 341 Cao, X., et al. (2014). "Anatomical, physiological and transcriptional responses of two  
342 contrasting poplar genotypes to drought and re-watering." Physiol Plant 151(4): 480-494.
- 343 Cohen, D., et al. (2010). "Comparative transcriptomics of drought responses in *Populus*: a  
344 meta-analysis of genome-wide expression profiling in mature leaves and root apices across two  
345 genotypes." BMC genomics 11(1): 630.
- 346 Dash, M., et al. (2018). "Gene network analysis of poplar root transcriptome in response to  
347 drought stress identifies a PtaJAZ3PtaRAP2.6-centered hierarchical network." PLOS ONE  
348 13(12): e0208560.
- 349 Davis, S. and P. S. Meltzer (2007). "GEOquery: a bridge between the Gene Expression  
350 Omnibus (GEO) and BioConductor." Bioinformatics 23(14): 1846-1847.
- 351 Fu, S. and M. C. Desmarais (2010). Markov blanket-based feature selection: a review of past  
352 decade. Proceedings of the world congress on engineering, Newswood Ltd.
- 353 Gillies, D. and A. Sudbury (2013). "Should causal models always be Markovian? The case of  
354 multi-causal forks in medicine." European Journal for Philosophy of Science 3(3): 275-308.
- 355 Grönlund, A., et al. (2009). "Modular gene expression in Poplar: a multilayer network  
356 approach." New Phytologist 181(2): 315-322.
- 357 Hamanishi, E. T., et al. (2015). "Poplar trees reconfigure the transcriptome and metabolome in  
358 response to drought in a genotype- and time-of-day-dependent manner." BMC genomics 16(1):  
359 329.
- 360 Han, X., et al. (2020). "Comparative transcriptome analyses define genes and gene modules  
361 differing between two *Populus* genotypes with contrasting stem growth rates." Biotechnology  
362 for biofuels 13(1): 1-21.
- 363 Jia, J., et al. (2016). "Physiological and transcriptional regulation in poplar roots and leaves  
364 during acclimation to high temperature and drought." Physiol Plant 157(1): 38-53.
- 365 Lahiri, A., et al. (2019). "Bayesian modeling of plant drought resistance pathway." BMC Plant  
366 Biology 19(1): 96.

367 Liber, Y., et al. (2020). "A Bayesian network approach for the identification of relationships  
368 between drivers of chlordecone bioaccumulation in plants." *Environmental Science and  
369 Pollution Research*: 1-6.

370 Liu, Q., et al. (2016). "PoplarGene: poplar gene network and resource for mining functional  
371 information for genes from woody plants." *Scientific Reports* 6(1): 31356.

372 Lorenz, W. W., et al. (2011). "Microarray analysis and scale-free gene networks identify  
373 candidate regulators in drought-stressed roots of loblolly pine (*P. taeda* L.)." *BMC genomics*  
374 12: 264.

375 Molina, C., et al. (2008). "SuperSAGE: the drought stress-responsive transcriptome of  
376 chickpea roots." *BMC genomics* 9(1): 553.

377 Ogata, Y., et al. (2009). "A database for poplar gene co-expression analysis for systematic  
378 understanding of biological processes, including stress responses." *Journal of Wood Science*  
379 55(6): 395-400.

380 Scutari, M., et al. (2014). "Multiple quantitative trait analysis using Bayesian networks."  
381 *Genetics* 198(1): 129-137.

382 Shen, J., et al. (2008). *Markov Blanket Feature Selection for Support Vector Machines*. AAAI.

383 Stolf-Moreira, R., et al. (2011). "Transcriptional Profiles of Roots of Different Soybean  
384 Genotypes Subjected to Drought Stress." *Plant Molecular Biology Reporter* 29(1): 19-34.

385 Street, N. R., et al. (2006). "The genetics and genomics of the drought response in *Populus*."  
386 *The Plant Journal* 48(3): 321-341.

387 Tan, Y. and Z. Liu (2013). *Feature selection and prediction with a Markov blanket structure  
388 learning algorithm*. BMC bioinformatics, Springer.

389 Tasaki, S., et al. (2015). "Bayesian network reconstruction using systems genetics data:  
390 comparison of MCMC methods." *Genetics* 199(4): 973-989.

391 Vignes, M., et al. (2011). "Gene regulatory network reconstruction using Bayesian networks,  
392 the Dantzig Selector, the Lasso and their meta-analysis." *PloS one* 6(12): e29165.

393 Wang, L., et al. (2019). "High-dimensional Bayesian network inference from systems genetics  
394 data using genetic node ordering." *Frontiers in genetics* 10: 1196.

395 Zhang, H. and T. Yin (2016). "Identifying candidate genes for wood formation in poplar based  
396 on microarray network analysis and graph theory." *Tree Genetics & Genomes* 12(3): 61.

397 Zhu, J., et al. (2007). "Increasing the power to detect causal associations by combining  
398 genotypic and expression data in segregating populations." *PLoS Comput Biol* 3(4): e69.

399

## Figures

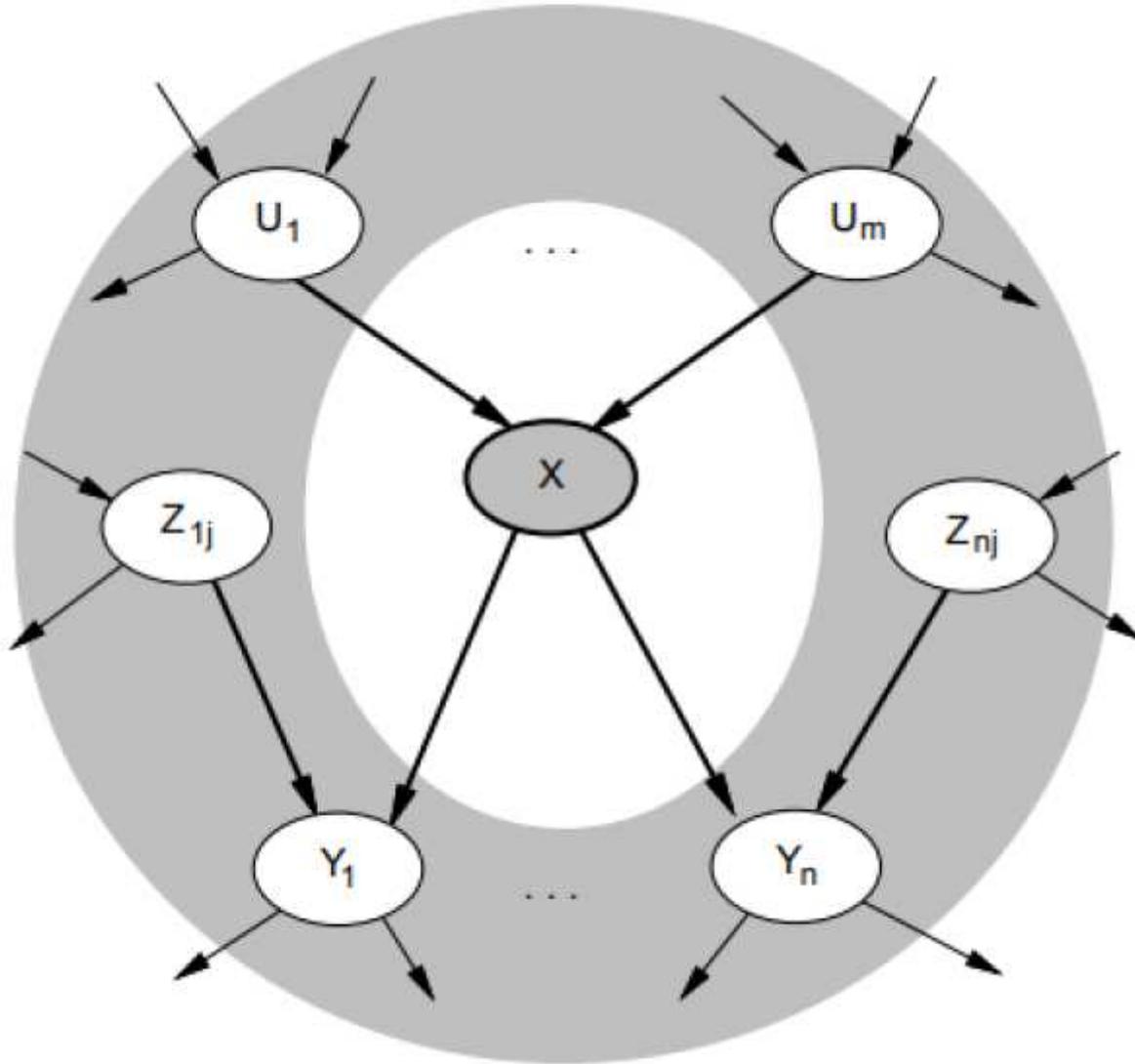
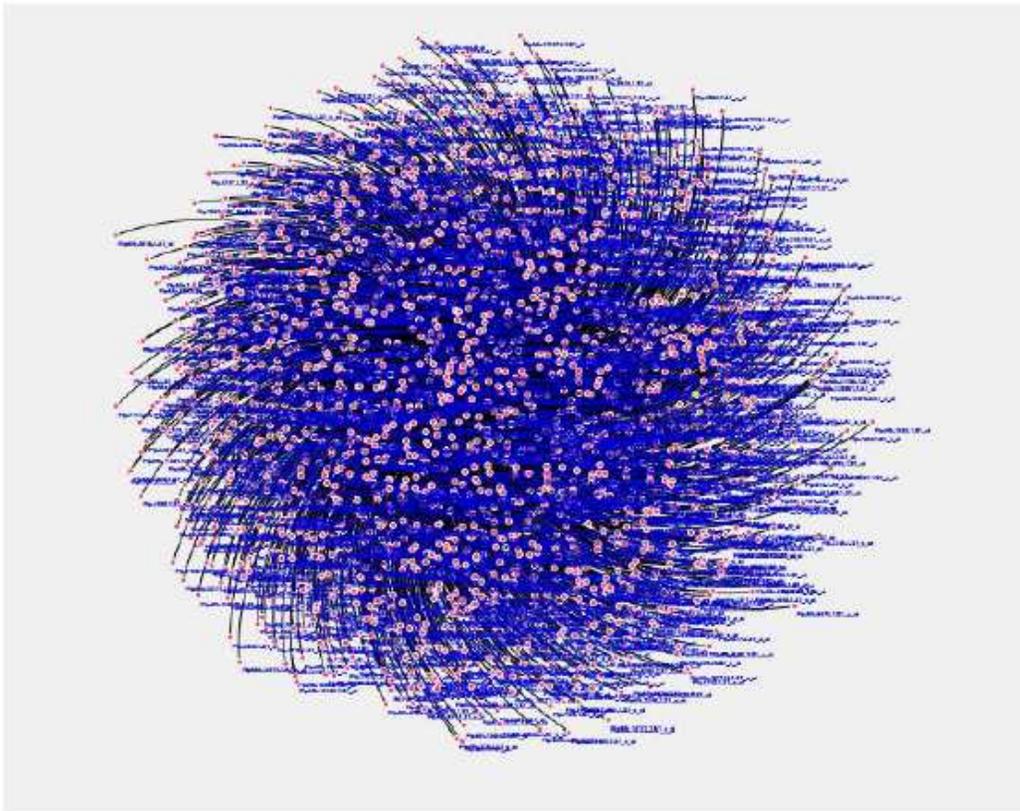
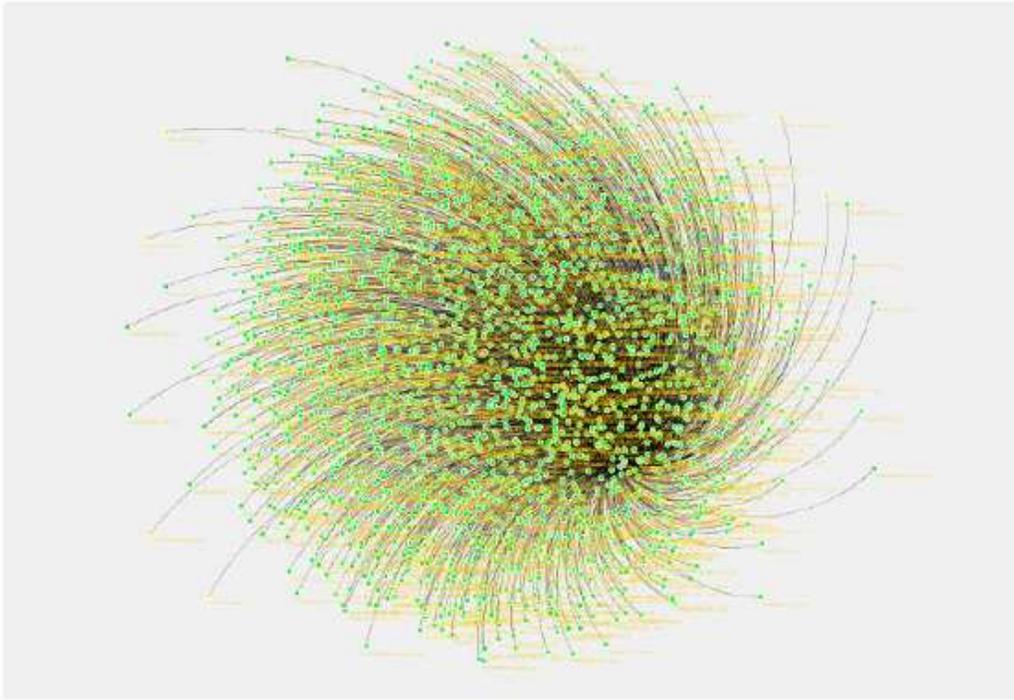


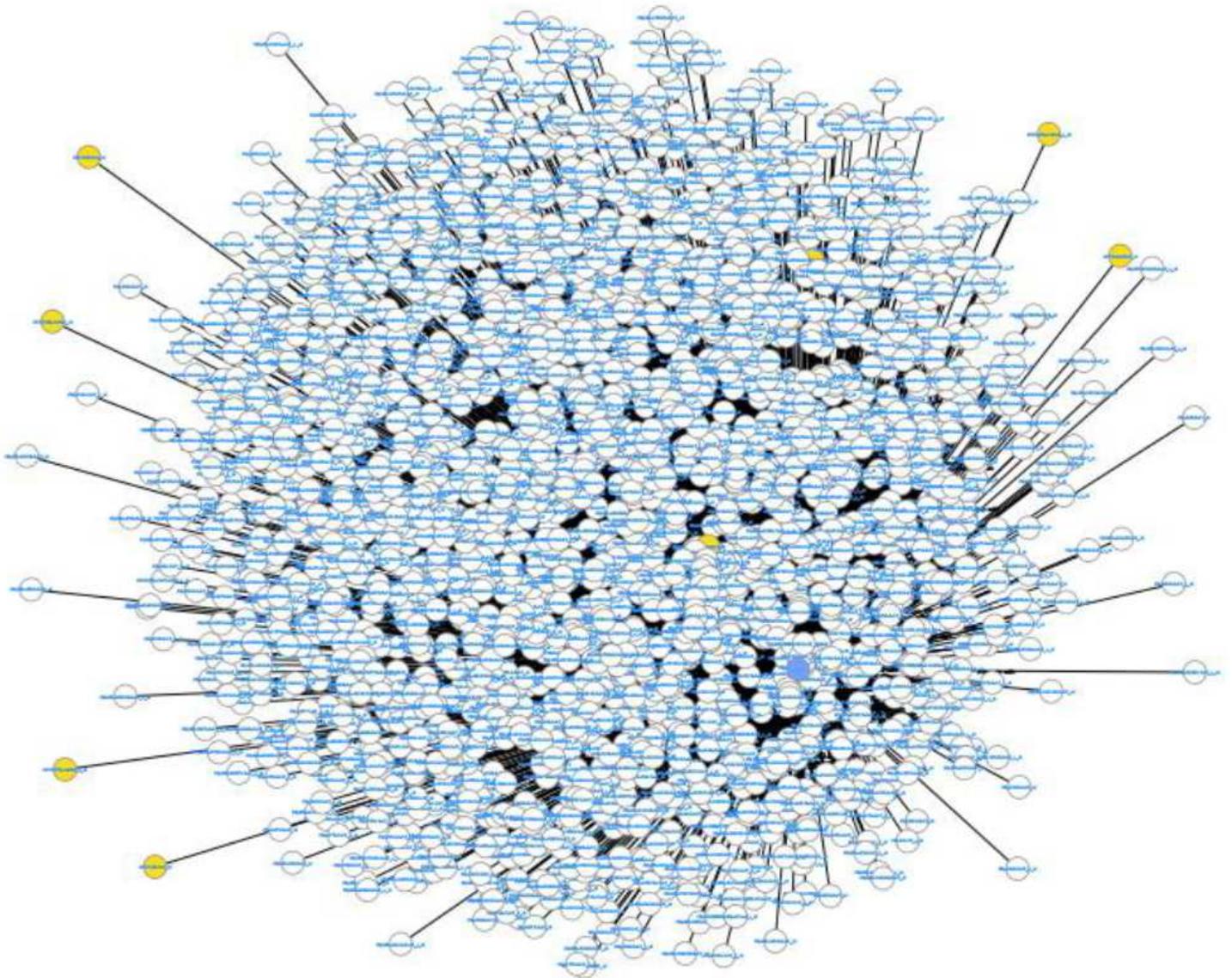
Figure 1

MB of a given postulated gene BN shows that MB of a node/probe/gene is the set containing the node/gene's parents, children, and co-parents.



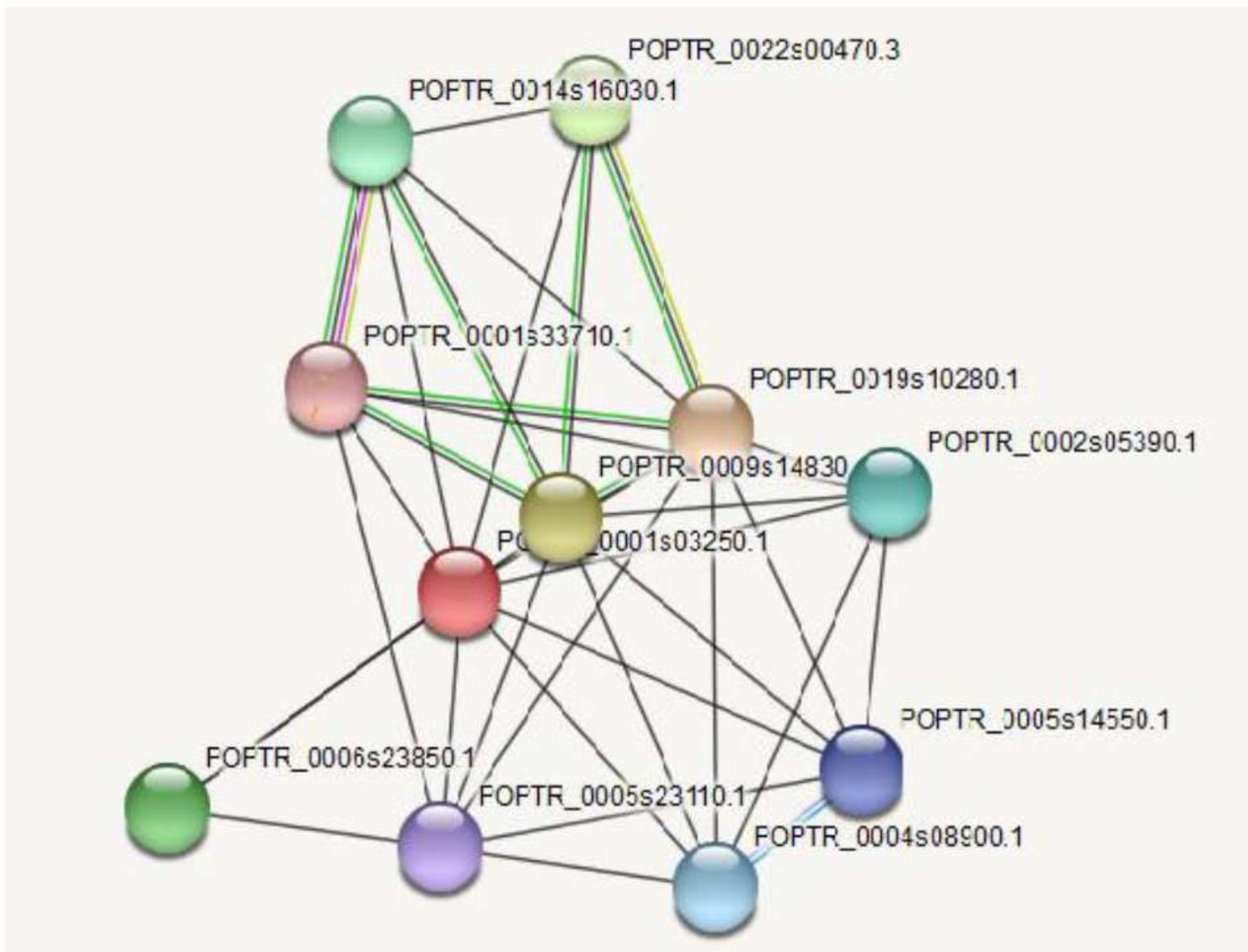
**Figure 2**

The snapshot of gene connectivity in learned BN in two graph forms.



**Figure 3**

Highlighting top 10 genes with their neighbors using MCC algorithm



**Figure 4**

Network of protein-protein interactions by STRING database information.

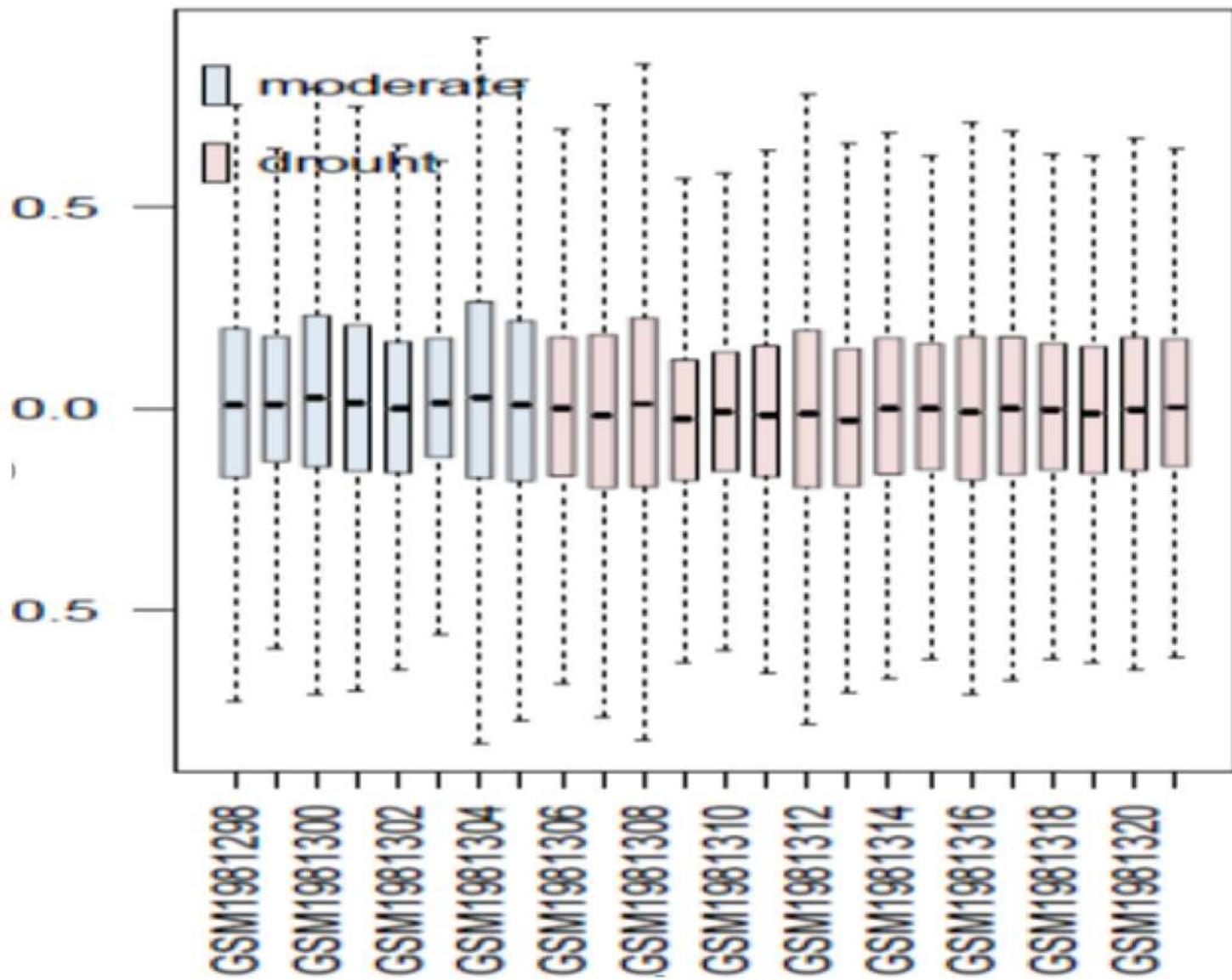


Figure 5

Boxplot of GSE764322 using in this study.