# GENE CLUSTERING AND CONSTRUCTION OF INTRA-CLUSTER GENE REGULATORY NETWORK

Aparajita Khan

REGISTRATION No.: 121256 OF 2012-13

EXAMINATION ROLL No.: M6TCT1518

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

JADAVPUR UNIVERSITY



#### Introduction

• Identification of Differentially Expressed Genes

- Clustering of Significant Genes
- Construction of Gene Regulatory Network
- Conclusion
- References





# **Genes and Gene Expression**

- DNA genetic material that stores and transmits hereditary information from generation to genera
- Gene segment of DNA that codes for a specific protein
- All cells in an organism have identical genes but gene expression and regulation is responsible for robust dynamic behavior of cell





## **Microarray Technology**

Monitors genome wide expression
 profile of organism under study

 Microarray – glass slide with DNA molecules orderly fixed at specific spots

 Spots contain few million copies of DNA corresponding to particular gene

• Works by DNA Hybridization



# Motivation of work

•Lung cancer - uncontrolled growth of abnormal cells in one or both lungs.

•Cancer develops due to DNA damage and disruption of control of gene expression and regulation.

•Gene expression analysis provides insight into gene function and disease physiology.

In this work 3 key issues are addressed

- identification of genes exhibiting differential gene expression patterns between diseased and non-diseased population
- clustering significant genes into groups of coexpressed genes
- construction intra-cluster gene- regulatory network

# Dataset Under Study

• The study [1] explores patterns of pathway deregulation in normal airway epithelial cells from patients with and without lung cancer.

• Gene Expression Omnibus (GEO) Dataset GDS2771 [2] reports genome wide expression values for 192 smokers with suspect of lung cancer.

- 2 groups of population in dataset
- <sup>•</sup>Healthy smokers (90 samples)
- Diseased smokers (97 samples)



# **Gene Subset Selection**

 COSMIC, the Catalogue Of Somatic Mutations In Cancer
 (<u>http://cancer.sanger.ac.uk</u>) manually curates 2002811 coding point mutations in human genes from scientific literature.

•COSMIC's *Cancer Browser* tool reports 24283 genes including alias to be mutated in lung cancer.

• Gene expression profile corresponding to 11237 genes were found in the dataset GDS2771 were considered for further study.

# IDENTIFICATION OF COMMONS DIFFERENTIALLY EXPRESSED GENES



•Differentially expressed genes - gene data determined to be statistical outliers from some standard state, which cannot be ascribed to chance or natural variability.

•Identify genes exhibiting differential expression patterns amongst the two classes of population

<sup>•</sup>These genes serve as potential pharmaceutical targets and diagnostic markers.



•**Null Hypothesis**  $H_0$ : expression levels in healthy and diseased population comes from normal distributions with equal means.

• Alternative hypothesis  $H_A$ : data comes from different distributions with unequal means.

• For each gene compute Welch's *t*-statistic  $t = \frac{\overline{X_A} - \overline{X_B}}{\sqrt{\frac{s_A^2}{N_A} + \frac{s_B^2}{N_B}}}$  where  $\overline{X_A}, \overline{X_B}$  are the sample means, and sample variances are  $s_A^2, s_B^2$ 

• p-value quantifies the probability of observing a test statistic of this extreme or more given that both samples come from the same distribution.

<sup>•</sup>A small p-value indicates that there is a small chance of getting this data if no real difference existed.

# **False Discovery Rate (FDR) and q-Value**

<sup>•</sup>p-value measures significance in terms of **false positive rate**.

• p-value cut-off of 0.05 means that there is 5% chance that this gene is called significant when it is truly null. Out of 15000 genes 750 are false positives.

•q-values measures significance in terms of false discovery rate.

$$FDR = E\left[\frac{F}{T+F}\right] = E\left[\frac{F}{S}\right]$$

<sup>•</sup>Takes into account that several features are simultaneously tested, hence a better measure.

<sup>•</sup>A q-value cut off of 0.05 results in a FDR of 5% among all significant genes.



Plot of number of significant genes filtered for respective p-value and q-value cutoff



### • Number of genes filtered for 4 statistically significant p and q value cutoffs

Cut off value	No. of Significant genes based on p-value	No. of Significant genes based on q-value
0.001	374	31
0.005	968	168
0.01	1231	316
0.05	2659	1270

<sup>•</sup>We use q-value cutoff of 0.005 and the resultant list of 168 genes is identified as differentially expressed between the diseased and healthy population.





•Clustering - statistical technique used to generate a category structure that fits a set of observations.

•High degree of association between members of the same group, low degree of association between members of different groups.

Clustering of gene expression patterns is used to identify groups of co-expressed genes.
 Genes belonging to same cluster are typically involved in related functions and are frequently co-regulated..



•Hierarchical clustering: agglomerative

**Distance** Measures

• Euclidean Distance

• Manhattan Distance  $D_{Man}(X,Y) = \sum_{i=1}^{N} |x_i - y_i|$ 

$$D_{Euc}(X,Y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$
$$D_{Euc}(X,Y) = \sum_{i=1}^{n} |x_i - y_i|$$

• Mahalanobis distance  $Dis_{Mah} = \sqrt{(X-Y)C^{-1}(X-Y)^T}$  where  $C_{ii} = cov(X_i, X_j) = E[(X_i - E[X_i])(X_j - E[X_j])]$  $Dis_{Corr} = 1 - \frac{\left(X - \overline{X}\right)\left(Y - \overline{Y}\right)}{\sqrt{\left(X - \overline{X}\right)\left(X - \overline{X}\right)^{T}}}\sqrt{\left(Y - \overline{Y}\right)\left(Y - \overline{Y}\right)^{T}}$ •Pearson correlation coefficient  $Dis_{Spr} = 1 - \frac{\left(R_X - \overline{R_X}\right)\left(R_Y - \overline{R_Y}\right)'}{\sqrt{\left(R_X - \overline{R_X}\right)\left(R_X - \overline{R_X}\right)^T}} \sqrt{\left(R_Y - \overline{R_Y}\right)\left(R_Y - \overline{R_Y}\right)^T}$ Spearman Rank correlation coefficient



# Linkage criterion and Dendrogram



5

4

3

2

5 1 3 2

•Cutting the dendrogram at different depths produces different number of clusters



• Hierarchical Clustering

Cophenetic Correlation Coefficient (CPCC) measures original pairwise dissimilarities between the feature vectors and the cophenetic

dissimilarities from the dendrogram

#### **Partitional Clustering**

- Ounn's Index
- Davies-Bouldin Index
- Calinski- Harabasz Index
- Silhouette Index
- CS Index
- SYM-Index
- SV-Index



# Distance and Linkage criterion Selection

#### Hierarchical Cluster Analysis on Healthy Population

	Single Linkage	Average Linkage	Complete Linkage
Euclidean Distance	0.7024	0. 8263	0.7483
Manhattan Distance	0.7022	0.8457	0.7586
Mahalanobis Distance	0.8205	0.8437	0.3039
Pearson Correlation	0.7071	0.8315	0.8128
Distance			
Spearman Correlation	0.6750	0.8484	0.7612
Distance			

#### **Hierarchical Cluster Analysis on Diseased Population**

	Single Linkage	Average Linkage	Complete Linkage
Euclidean Distance	0.6515	0.8254	0.7230
Manhattan Distance	0.6669	0.8172	0.7182
Mahalanobis Distance	0.8319	0.8550	0.3216
Pearson Correlation	0.8844	0.9118	0.8653
Distance			
Spearman Correlation	0.9046	0.9224	0.8843
Distance			



#### For Healthy Population

k-parameter	Dunn's Index	Davies- Bouldin Index	Calinski- Harabasz Index	Silhouette Index	CS Index	SYM Index	SV Index
2	0.519489466	0.421794	179.4794	0.320344	2.417697	0.123279	0.058357
3	0.521822826	1.285558	240.179	0.341226	3.007731	0.202354	0.475018
4	0.301566338	1.354472	127.1219	0.267722	3.691343	0.083151	0.375878
5	0.301566338	1.235196	96.22031	0.249841	2.543032	0.11748	0.397804
6	0.301566338	1.175929	77.90129	0.233563	2.130993	0.187262	0.417877
7	0.301566338	1.100681	65.56314	0.226116	2.929643	0.169686	0.275173
8	0.301566338	1.026924	56.36059	0.214371	4.345251	0.150266	0.187871
9	0.250779006	1.204307	51.56966	0.22683	5.736134	0.133218	0.151053
10	0.250779006	1.215522	46.92917	0.225899	6.390402	0.122311	0.132319



### For Diseased Population

k-parameter	Dunn's Index	Davies- Bouldin Index	Calinski- Harabasz Index	Silhouette Index	CS Index	SYM Index	SV Index
2	0.323957822	0.546018	102.3888	0.3014	2.51132	0.0931	0.093128
3	0.339238222	0.900783	137.9075	0.3158	2.472717	0.1532	0.153163
4	0.375366245	1.121942	240.9651	0.3796	1.976578	0.2221	0.122059
5	0.339434355	1.061795	80.1508	0.2891	5.670765	0.0992	0.099245
6	0.296278403	1.031703	64.8691	0.2763	6.60191	0.1524	0.152398
7	0.296278403	0.965214	54.5249	0.2595	1.219989	0.1577	0.157732
8	0.296278403	0.883861	47.1872	0.2603	1.067121	0.1773	0.177312
9	0.352699333	0.936416	44.0637	0.2698	9.959585	0.1571	0.157077
10	0.352699333	0.850708	39.3389	0.274	8.445925	0.147	0.147001

# - Dendrogram for optimal k-parameter

For Healthy Population

**For Diseased Population** 



Cluster Assignment Analysis		
<sup>•</sup> For healthy population genes distributed as	<sup>•</sup> For diseased population genes of	listributed as
Cluster#1 20 genes	Cluster#1 13 genes	Cluster#3 60 genes
Cluster#2 72 genes	Cluster#2 18 genes	Cluster#4 77 genes

Cluster#3 76 genes

<sup>•</sup>Cluster# 1 of Healthy Dataset and Cluster# 2 of Diseased Dataset 72.22% similarity

<sup>•</sup>Cluster# 2 of Healthy Dataset and Cluster# 3 of Diseased Dataset 100% similarity

<sup>•</sup>Cluster# 3 of Healthy Dataset and Cluster# 4 of Diseased Dataset 89.61% similarity

<sup>•</sup>Extra cluster of diseased population: 10 genes (Cluster# 2), 2 genes (Cluster# 1) and 1 gene (Cluster# 3) of healthy dataset

<sup>•</sup>We identify 12 genes showing different cluster assignments in healthy and diseased population.



Gene Ontology

<sup>•</sup>Provides descriptions about gene products in terms of

biological process (BP)

molecular function (MF)

cellular components (CC)

- GO Term Finder tool (<u>http://go.princeton.edu/</u>)
- Significant GO terms shared between queried genes
- p-value : degree of enrichment
- Small p-value stronger evidence of annotation



# Healthy Population Cluster Analysis Cluster#1 : regulatory cluster

GOID	GO TERM from biological_process Ontology	p-value < 0.05	% of Genes of Cluster Annotated
GO:0044237	cellular metabolic process	0.004534	75
GO:0051128	regulation of cellular component organization	0.000635	30
GO:0048583	regulation of response to stimulus	0.014725	30
GO:0022414	reproductive process	0.000166	25
GO:0009653	anatomical structure morphogenesis	0.008663	25
GO:0080134	regulation of response to stress	0.002823	25
GO:0009719	response to endogenous stimulus	0.005762	20
GO:0006357	regulation of transcription from RNA polymerase II promoter	0.019413	20
GO:0002697	regulation of immune effector process	0.000658	15
GO:0050679	positive regulation of epithelial cell proliferation	0.001952	10



# Healthy Population Cluster Analysis Cluster#2 : response and signaling cluster

GOID	GO TERM from biological_process	p-value <0.001	% of Genes of Cluster Annotated
GO:0050896	response to stimulus	1.42E-06	65.55556
GO:0007154	cell communication	8.44E-05	53.05556
GO:0044700	single organism signaling	0.00014	41.66667
GO:0007165	signal transduction	0.000218	38.88889
GO:0044707	single-multicellular organism process	1.64E-05	37.5
GO:0048518	positive regulation of biological process	5.19E-06	36.11111
GO:0010468	regulation of gene expression	0.000363	31.94444
GO:0044767	single-organism developmental process	0.00053	30.55556
GO:0006950	response to stress	4.40E-05	29.16667
GO:0002376	immune system process	2.78E-05	22.22222
GO:0010646	regulation of cell communication	0.0004	22.22222
GO:0035556	intracellular signal transduction	0.000794	22.22222
GO:0006955	immune response	2.41E-05	18.05556
GO:0006952	defense response	0.000126	16.66667



# Healthy Population Cluster Analysis

## •Cluster#3 : cell development and maintenance cluster

GOID	GO TERM from biological_process Ontology	p-value <0.001	% of Genes of Cluster Annotated
GO:0019222	regulation of metabolic process	0.000159	43.42105
GO:0044767	single-organism developmental process	6.34E-06	35.52632
GO:0048856	anatomical structure development	6.93E-06	32.89474
GO:0030154	cell differentiation	0.000472	32.36842
GO:0048513	organ development	6.15E-05	21.05263
GO:0008283	cell proliferation	0.000142	15.78947
GO:0008219	cell death	0.000648	15.78947
GO:0010941	regulation of cell death	0.00099	13.15789
GO:0043687	post-translational protein modification	2.64E-06	9.210526
GO:0006281	DNA repair	0.000851	9.210526
GO:0044770	cell cycle phase transition	3.33E-05	9.210526
GO:0010564	regulation of cell cycle process	6.58E-05	9.210526
GO:1901214	regulation of neuron death	0.000786	5.263158



# **GENE REGULATORY NETWORK**



• GRN - most important organizational level in the cell, signals from the cell state and the outside environment are integrated in terms of activation and inhibition of genes.

• "Connection" connotes regulatory interaction

•Genes having similar gene expression profiles are more likely to regulate one another or be regulated by some other common parent gene.

<sup>•</sup>Bayesian Network approach

Modified version of Sparse Candidate Algorithm [6]

# – Bayesian Network Formalisms

• Represents joint probability distribution of random variables consisting of two components  $B = \langle \mathbf{G}, \boldsymbol{\varphi} \rangle$ 

• **G** -DAG vertices correspond to  $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$ 

• Encodes Markov assumption  $\forall X_i (X_i \perp Non - Descendants_{X_i} | \mathbf{U}_{X_i})$ 

•  $\boldsymbol{\varphi}$  conditional probability distribution (CPD) of  $P(X_i | \mathbf{U}_{X_i})$ 

• Joint distribution satisfying conditional independence properties decomposed by chain rule for Bayesian networks  $P(X_1, X_2, \dots, X_N) = \prod_{i=1}^{N} P(X_i | \mathbf{U}_{X_i})$ 



# – Learning Bayesian Network

• Given a training set of samples  $D = x[1], x[2], \dots, x[M]$  independently drawn from some unknown Bayesian Network  $G^*$  with underlying distribution  $P^*$ , goal : recover  $G^*$ 

• Parameter Learning : given network structure G, set of data instances D determine what values of  $\varphi$  best describes the generated data.

• Maximum Likelihood estimation : 
$$L(\varphi:D) = \prod_{m=1}^{M} P(\mathbf{x}[m] | \varphi)$$
, choose  $\overline{\varphi} = \max_{\varphi} L(\varphi:D)$  closed form :  $\overline{\varphi} = \frac{M[x, u]}{M[u]}$ 

• **Structure Learning** :define scoring function that measures how well each model fits the observed data. Optimization algorithm employed to search for highest scoring model.

• By Bayes rule,  $P(\mathbf{G} | \mathbf{D}) = \frac{P(\mathbf{D}/\mathbf{G})P(\mathbf{G})}{P(\mathbf{D})}$  ignoring normalization factor  $score_B(\mathbf{G}:\mathbf{D}) = \log P(\mathbf{D}/\mathbf{G}) + \log P(\mathbf{G})$ where  $P(\mathbf{D}/\mathbf{G}) = \int P(\mathbf{D}/\mathbf{G}, \boldsymbol{\varphi})P(\boldsymbol{\varphi}/\mathbf{G}) d\boldsymbol{\varphi}$  and by decomposability  $score_B(\mathbf{G}:\mathbf{D}) = \sum_{i=1}^{N} Score_Contribution(X_i, \mathbf{U}_{X_i}^{\mathbf{G}}:\mathbf{D})$ 

# - Search Algorithm

•**Input** : training set , scoring function, set of possible network structures

•**Output** : network structure that maximizes score

•  $2^{o(n^2)}$  possible structures : NP-Hard problem

<sup>•</sup>Initial Network Formation

• Discretization : **Z-score** 

• Mutual Information 
$$\operatorname{MI}(X_i, X_j) = \sum_{x,y} \hat{P}(x, y) \log \frac{\hat{P}(x, y)}{\hat{P}(x)\hat{P}(y)}$$

• Select k nodes giving highest MI as candidate parents for  $X_i$ 

# Cycle\_Removal(Bn) Repeat until $B_n$ is not a DAG Detect cycle in graph using DFS : $C = X_n \rightarrow X_{i2} \rightarrow \cdots \rightarrow X_{im} \rightarrow X_n$ List nodes involved in cycle $L = \{X_{i1}, X_{i2}, \dots, X_{in}\}$ $\forall X_{ii} \in L \text{ find } parent(X_{ii}) \text{ in } C$ Find rank of $parent(X_{ii})$ as a candidate parent of $X_{ii}$ Construct $R = \{r_{i1}, r_{i2}, \cdots, r_{in} : r_{ij} = rank(parent(X_{ij})) \text{ in } C\}$ Find $r_{in} = \max\{R\}$ and $X_t = parent(X_{in})$ in C Remove edge $X_t \to X_{ij}$ and update $B_n = B_n - edge(X_t \to X_{ij})$

Return B<sub>n</sub>

-

# **Sparse Candidate Algorithm**

Loop for  $n = 1, 2, \cdots$  until convergence

#### Restrict:

#### Input:

- The training Dataset  $D = \{x[1], x[2], \dots, x[M]\}$
- An initial network  $B_{\theta}$

• A decomposable score such that  $score(B | D) = \sum_{i=1}^{N} Score \_Contribution(X_i, U_{X_i}^B : D)$ 

- A parameter k = maximum parent for each node
- A parameter *l* = minimum number of edges in network

Output: A network B

Based on D and  $B_{n-1}$  for each variable select a set  $C_i^n$  of candidate parents with  $|C_i^n| \leq k$  by *Candidate\_Parent\_Selection*  $(X_i, B_{n-1}, D, k)$ This defines a directed graph  $H_n = (\mathbf{X}, E)$  where  $E = \{X_j \rightarrow X_i | \forall i, j X_j \in C_i^n\}$ Remove cycles in  $H_n$  using *Cycle\_Removal*  $(H_n)$ 

#### Maximize:

Find network  $B_n = \langle \mathbf{G}_n, \boldsymbol{\varphi} \rangle$  maximizing  $score(B_n | D)$  among networks satisfying  $\mathbf{G}_n \subset H_n$  (i.e.  $\forall X_i, U_{X_i}^{\mathbf{G}_n} \subseteq C_i^n$ ) using *Greedy\_Hill\_Climbing*. Restrict the minimum number of edges by l to avoid sparse graph



Initial network based on mutual information between genes, 60 edges



Final GRN for Cluster# 1 of healthy population, 39 edges





Gene Interaction Network returned by GeneMANIA

Serial #	GeneMANIA Network Edge or Path	Network Study for the edge	Bayesian Con	Network Edge nection	Connectivity between nodes in Bayesian network	Dependence Evident From
1	C6 – NELL2	Ramaswamy- Golub-2001	$C6 \rightarrow$	NELL2	Parent-Child	Causal Reasoning
2	C6 – HTRA1	Burington- Shaughnessy- 2008	$C6 \rightarrow NEL$	$L2 \rightarrow HTRA1$	Descendant	Causal flow of dependence
3	BICC1-SOX9	Burington- Shaughnessy- 2008; Innocenti- Brown-2011	C Z SOX9	C6 S BICC1	Common Parent	Evidential reasoning causal reasoning
4	BICC1 –SOX9 – NR2F1	Kang-Willman- 2010	BICC1-	→ NR2F1	Parent-Child	Causal Reasoning
5	CLGN –ODF2	Mallon-McKay- 2013	DNA ØDF2	AJC6 S CLGN	Common Parent	Evidential reasoning causal reasoning



•Genome wide expression profile analysis of healthy and diseased population

 Identification of genes exhibiting differential expression pattern between healthy and diseased population

•Clustering genes into groups of co-expressed genes

•Construction of intra-cluster gene regulatory network

#### **Future Scope**

- Incorporation of fuzzy clustering measures
- Model GRN using continuous variables
   Bayesian network
- Inter-Cluster GRN with efficient learning algorithms
- Analysis of network connectivity change between healthy and diseased population



[1] Spira A, Beane JE, Shah V, Steiling K et al. "Airway epithelial gene expression in the diagnostic evaluation of smokers with suspect lung cancer". Nat Med 2007 Mar;13(3). [PMID: <u>17334370</u>]

[2] http://www.ncbi.nlm.nih.gov/sites/GDSbrowser?acc=GDS2771

[3] Storey, J.D., and Tibshirani, R.(2003), "Statistical significance for genomewide studies", PNAS 100(16), pp 9440–9445.
[4] Daphne Koller and Nir Friedman, *Probabilistic Graphical Models : Principles and Techniques*. MIT Press, 2009.
[5] Nir Friedman, Michal Linial, Iftach Nachman and Dana Pe'er, "Using Bayesian Networks to Analyze Expression Data", Journal of Computational Biology. August 2000, 7(3-4), pp 601-620.

[6] N. Friedman, I. Nachman, and D. Pe'er. "Learning Bayesian network structure from massive datasets: The "sparse candidate" algorithm". In Proceedings of the Fifteenth Annual Conference on Uncertainty in Artificial Intelligence (UAI), pp 196–205. 1999.

